

## Knowledge-based Avoidance of Drug-Resistant HIV Mutants

Richard H. Lathrop, Nicholas R. Steffen, Miriam P. Raphael,  
Sophia Deeds-Rubin, Michael J. Pazzani

Department of Information and Computer Science,  
University of California, Irvine, CA 92717 USA  
{ rickl | nsteffen | mraphael | sophiadr | pazzani }@uci.edu

**Paul J. Cimoch**

Director, Center for Special Immunology  
100 Pacifica #100, Irvine, CA 92718 USA  
gr82bpnc@msn.com

**Darryl M. See, Jeremiah G. Tilles**

Department of Medicine,  
University of California, Irvine, CA 92717 USA  
{ dmsee | jgtilles }@uci.edu

### Abstract

We describe an artificial intelligence (AI) system (CTSHIV) that connects the scientific AIDS literature describing specific HIV drug resistances directly to the Customized Treatment Strategy of a specific HIV patient. Rules in the CTSHIV knowledge base encode knowledge about sequence mutations in the HIV genome that have been found to result in drug resistance in the HIV virus. Rules are applied to the actual HIV sequences of the virus strains infecting the specific patient undergoing clinical treatment in order to infer current drug resistance. A search through mutation sequence space identifies nearby drug resistant mutant strains that might arise. The possible drug treatment regimens currently approved by the US Food and Drug Administration (FDA) are considered and ranked by their estimated ability to avoid identified current and nearby drug resistant mutants. The highest-ranked treatments are recommended to the attending physician. The result is more precise treatment of individual HIV patients, and a decreased tendency to select for drug resistant genes in the global HIV gene pool. The application is currently in use in human clinical trials on HIV patients. Initial results from a small clinical trial are encouraging and further clinical trials are planned. From an AI viewpoint the case study demonstrates the extensibility of knowledge-based systems because it illustrates how existing encoded knowledge can be used to support new applications that were unanticipated when the original knowledge was encoded.

### Problem Description

Human immunodeficiency virus (HIV) causes progressive deterioration of the immune system leading almost invariably to AIDS and death from opportunistic cancers and infections. Currently in the USA it is estimated to infect 3–5 million persons, is the leading cause of death in adults from age 14 to 35, and is the nation's leading cause of productive years of life lost aggregated over all age groups. HIV is estimated to infect 40–50 million persons worldwide (CDC 1997).

The high rate of HIV viral mutation both makes development of a vaccine difficult and results in rapid positive selection for drug resistant mutant strains. Recent multi-drug combination therapies are encouraging, but in most cases ultimately fail due to the development of drug resistance (O'Brian *et al.* 1996). A general theory of HIV drug resistance still is not in hand, but a number of specific sequence mutations in the HIV genome have been described in the scientific literature and associated with increased resistance to certain drugs.

In this paper we describe an AI system intended to improve the clinical treatment of individual HIV patients by identifying drug resistance in advance and avoiding it in treatment. This is done by first identifying drug resistant HIV mutant strains that already exist in the patient or are likely to be positively selected for by certain treatments, and then recommending a customized treatment designed to avoid selection of such mutants. The result is more precise treatment of individual HIV patients, and a decreased tendency to select for drug resistant genes in the global HIV gene pool.

## Project Goals

The project goals are:

1. Connect knowledge contained in the scientific literature about HIV drug resistance directly to the treatment of individual HIV patients;
2. Enable customized treatment strategies to be based on the HIV genotype that currently infects an individual HIV patient;
3. Identify the nature and extent of drug resistance currently present in an individual HIV patient;
4. Identify nearby drug resistant mutant strains that could be positively selected for by some treatments;
5. Rank the possible FDA-approved treatments by an estimate of their ability to avoid both current and nearby drug resistant mutants;
6. Estimate the costs of the highest-ranked treatments;
7. Recommend treatments that are estimated to be most likely to avoid known HIV drug resistance.

## Related Work

An expert system based on experimental data from HIV patients (immunologic markers) has been used to diagnose the opportunistic non-Hodgkin's lymphomas which often develop (Diamond *et al.* 1994). Knowledge-based systems have been applied to HIV patient medical record systems (Musen *et al.* 1995; Safran *et al.* 1996), monitoring ongoing HIV patient protocols (Musen *et al.* 1996; Tu *et al.* 1995; Sonnenberg, Hagerty, & Kulikowski 1994; Sobesky *et al.* 1994), and HIV patient assessment (Xu 1996; Ohno-Machado *et al.* 1993). Less closely related are knowledge-based systems that apply qualitative modeling and process simulation to HIV laboratory systems (Sieburg 1994; Ruggiero *et al.* 1994). To our knowledge CTSHIV is the first system to use HIV sequence data from HIV patients to estimate current and nearby drug resistant mutants and recommend treatment combinations to avoid both.

## Domain Background

The information content of an HIV virus is contained in a set of genes encoded in its genome. Each gene is a sequence of bases or nucleotides of four varieties. A gene can be represented as a string over an alphabet of four characters, one character representing each nucleotide. The HIV genome ultimately causes the production of gene products, often proteins, important in the virus life cycle. A protein is a sequence of amino acids of twenty varieties, and can be represented as a string over an alphabet of twenty characters. Each amino acid in the protein is encoded by a block of three adjacent nucleotides in the genome, called a codon. The two proteins targeted by current FDA-approved drugs are called "reverse transcriptase" (RT) and "protease" (PRO). An example RT protein structure (Hsiou *et al.* to be published) is shown in Figure 1.

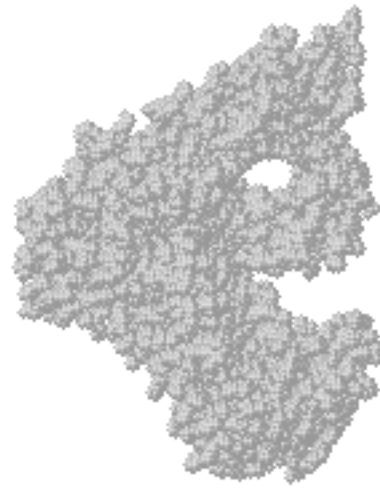


Figure 1: The 3D structure of the HIV reverse transcriptase protein (PDB code 1DLO). Each sphere represents an atom. The structure is encoded in the HIV reverse transcriptase sequence (see Figure 3). Mutations in the sequence cause changes in the number, type, or spatial arrangement of atoms in the structure.

The genome string must be copied from one generation to the next during the virus life cycle. Copying errors occur frequently, and are called mutations. Mutations can change the structure or function of the virus, and thus alter how it interacts with its environment. Mutant strains with genome sequences very similar to the patient's current strain (close in Hamming or edit distance) appear spontaneously and continuously. In a full-blown case of AIDS, it is estimated that every single point mutation appears every day, every coordinated pair of point mutations appears once or more during the course of the infection, and even coordinated triples of point mutations may appear (Condra *et al.* 1995).

A drug typically works by blocking a key part of the virus life cycle. A drug resistant mutation occurs when a copying error in the viral genome so alters the virus that it can perform the targeted step of its life cycle even in the presence of the drug. In the continued presence of the drug the mutant strain may out-compete the dominant strain, and thereby may itself become the dominant strain in the patient. This is often called selective drug resistance, because the resistant mutant is selected for by the drug's presence. If unrecognized, the current treatment may lose its effect and the patient's condition may deteriorate. The resulting strain is more challenging to treat because the treatment options have been reduced. If the drug treatment is changed in response, the potential for a new drug resistant mutation to develop is present. The use of an increasing variety of drugs has led to virus strains increasingly resistant to multiple drugs simultaneously. Sadly, increasing prevalence of drug resistant strains in the HIV global gene pool means that new patients may be infected by mu-

tant strains that already have accrued resistance from previous hosts (Gu *et al.* 1994). Consequently it is important to avoid selecting for drug resistant mutants.

Combination treatments involving multiple drugs are one approach to avoiding drug resistance (Lange 1995). If the virus mutates to resist one drug but still is inhibited by another, it may be suppressed or unviable. In this case the mutation may not be positively selected for. Combinations may contain up to four simultaneous drugs, but usually do not exceed three due to the potential for intolerable side-effects and toxicity. Severe side-effects often induce a patient to stop one or more drugs without knowledge of their physician, called non-adherence (formerly non-compliance). Non-adherence negates combination therapy and increases the likelihood of selecting for drug resistant mutants.

Combinations containing at least one protease inhibitor are referred to as Highly Active Anti-Retroviral Therapy (HAART). HAART typically results in a dramatic drop in viral load within two weeks, often sustained for long periods of time. Enthusiasm for the potential of HAART to eradicate HIV has been tempered by the inevitable failure of these regimens due to the eventual development of drug resistance (Carpenter *et al.* 1996). The virus appears to remain in a proviral state in resting memory T-cells, where it is inaccessible to antiretroviral drugs (Wong *et al.* 1997; Finzi *et al.* 1997). Mutations still can occur under HAART, though the mutation rate is greatly decreased (Jacobsen *et al.* 1996).

There are important limitations of the approach below. Sequence-based rules capture only part of the domain knowledge about drug resistance, albeit a clinically useful part. Drug resistance may arise for other domain-specific reasons that cannot be represented easily as rules. Current sequencing techniques may provide only partial or no information about minority strains. The rule set is only as complete as current scientific knowledge allows. Currently it may be possible to infer when resistance is likely to occur, based on genome sequences actually seen in the patient that correspond to resistance-conferring mutations described in the scientific literature. However, it is impossible to guarantee the non-existence of an unsuspected resistant mutant.

Nonetheless, knowledge of current or nearby mutants putatively resistant to one or more drugs is valuable to a physician treating an HIV patient. In conjunction with HAART, such knowledge may help select a combination of drugs less likely to be resisted. Currently there are 11 drugs approved by the FDA for HIV, plus one available for compassionate use. These 12 result in 407 different combination treatments of four or fewer drugs, as some drugs should not be used together. A physician may find it tedious to scan many sequences, be unfamiliar with the latest HIV drug resistant mutations reported, or have difficulty ranking the hundreds of treatment choices for each patient. CTSHIV mediates between the scientific literature and the patient's current infection to help a physician avoid HIV drug resistance.

## Application Description

The application (1) accepts as input experimentally determined HIV sequences extracted from the patient; (2) extracts the relevant codons and constructs virtual genomes; (3) estimates current resistance by applying knowledge base rules; (4) searches nearby mutation sequence space to identify nearby putatively resistant mutants; (5) ranks the possible FDA-approved treatment regimens according to their ability to avoid selective drug resistance; and (6) recommends the highest-ranked treatment regimens to the attending physician. See the application overview flowchart in Figure 2.

## Patient's Experimental Data

The reverse transcriptase and protease portions of the POL gene are amplified from each patient. Clones are produced, plasmid DNA is extracted, and the sequence is determined using a commercially available ABI sequencer. The reverse transcriptase sequence contains 1,299 letters (433 codons) and the protease sequence contains 297 letters (99 codons). Figure 3 shows an example HIV sequence from an HIV patient.

The sequences are pre-aligned to a standard reference HIV sequence (HXB2) using standard sequence alignment algorithms. Deviations from the reference sequence correspond to mutations in the virus infecting the patient. Typically five reverse transcriptase sequences and five protease sequences, a total of 7,980 letters of HIV genomic information, are the input experimental data on the patient's current infection.

## Extract Features, Objects

Processing in this step is routine. The features extracted are exactly those codons in positions referred to by the antecedent of some rule. Other positions are not yet associated with known drug resistance. Currently 55 rules mention 31 different codon positions, 20 in RT and 11 in PRO. HIV sequences are replaced by abstract objects consisting of only those codon positions. All possible virtual genomes are formed consistent with the experimental sequences.

## Identify Current Resistance

Current drug resistance is identified by applying the 55 rules in the knowledge base to the HIV sequences from the patient. The rules represent knowledge about HIV drug resistance as a set of if-then rules of the form:

IF  $\langle$  antecedent  $\rangle$  THEN  $\langle$  consequent  $\rangle$  [weight].

For example, one such rule in CTSHIV is:

**IF Methionine is encoded by RT codon 151,  
THEN do not use AZT, ddI, d4T, or ddC.  
[weight= 1.0]  
(Iversen *et al.* 1996)**

The weight associated with a rule is not a confidence as in many expert systems. Rather, it reflects the estimated level of resistance to a particular drug, and is part of the consequent. Weights range from 0.1 (low)

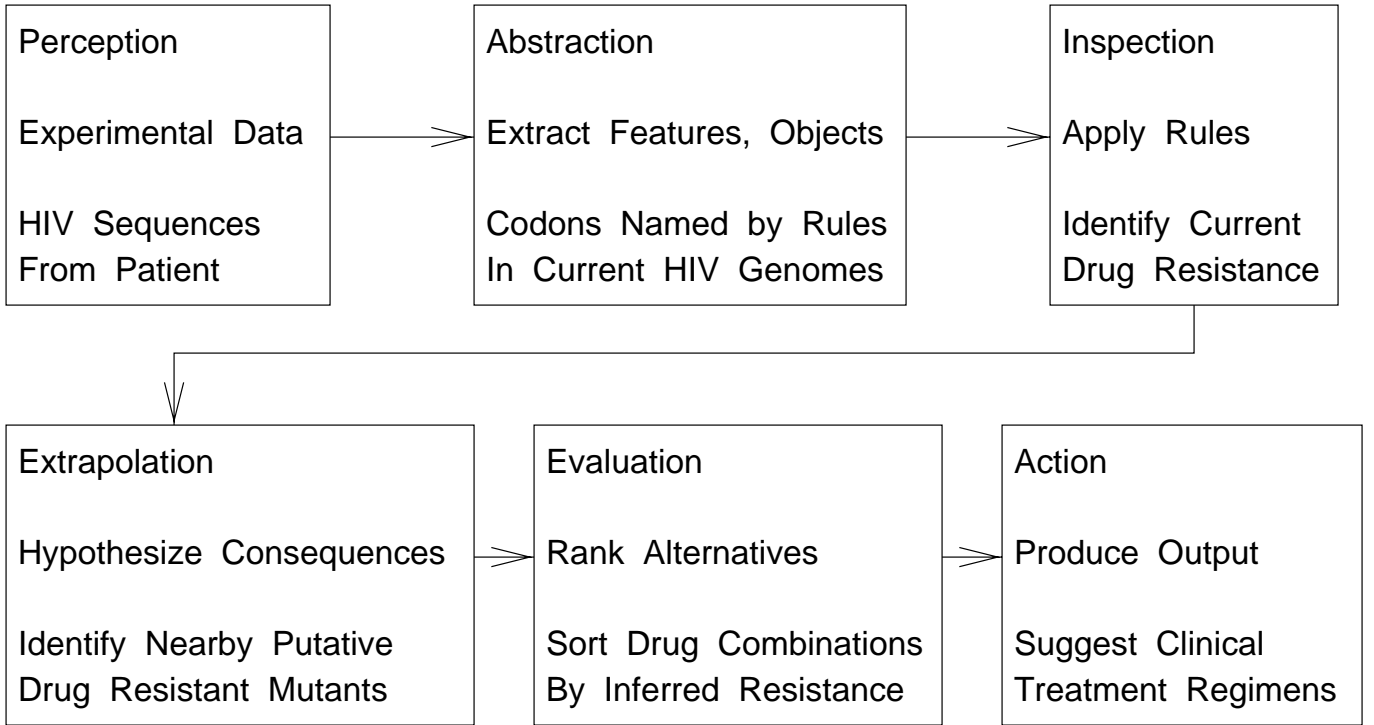


Figure 2: Application overview flowchart.

to 1.0 (high) based upon expert advice and the level of resistance reported in the literature.

To estimate current resistance, rule weight is multiplied by the fraction of viral sequences that trigger the rule, and combined additively. As a summary metric we use

$$CurrWt(D) = \sum_{r \in Rules(D)} \sum_{s \in S} Apply(r, s) / |S|$$

where  $D$  is a set of drugs comprising a combination therapy,  $Rules(D)$  are the rules that confer resistance to a drug in  $D$ ,  $S$  is the set of the HIV sequences extracted from the patient, and  $Apply(r, s)$  yields the rule weight of  $r$  if  $r$  fires on  $s$  and 0 if not.  $CurrWt$  is comparable only between combinations with the same number of drugs because any superset of a drug combination has equal or greater current weight.

Under this model, the total current level of resistance to a multi-drug combination is the sum of the current resistances to each drug. The effect of this is to identify drug combinations that have little or no current resistance and therefore attack the virus strongly.

### Predict Nearby Resistant Mutants

Nearby resistant mutants are predicted by a backward-chaining search through mutation sequence space, beginning with the patient's current HIV sequences. At each step, a sequence that does not fire a rule is used

to generate several new sequences that do. The new sequences are identical except that codon positions mentioned by the rule are modified so that the rule does fire. They represent mutants that are close in Hamming distance but resist the drugs mentioned by the rule. Conceptually, every virtual mutant within a pre-determined Hamming distance cut-off is examined. Currently all mutants up to and including Hamming distance three are considered. Branch-and-bound techniques speed the search by pruning unnecessary examinations. Currently CTSHIV runs in about a minute per patient, which is acceptable for now.

To predict nearby mutants, rule weights are combined by taking the maximum across all mutants of the minimum across all drugs in the combination. As a summary metric we use

$$\begin{aligned} \mathbf{m\_dist}(D) &= \min\{h | \exists x \in M(S, h), \forall d \in D, \\ & 0 < \sum_{r \in Rules(d)} Apply(r, x)\} \end{aligned}$$

$$\begin{aligned} \mathbf{m\_wt}(D) &= \max_{x \in M(S, \mathbf{m\_dist}(D))} \min_{d \in D} \sum_{r \in Rules(d)} Apply(r, x) \end{aligned}$$

$$\begin{aligned} MutScore(D) &= \max\{0, h_{max} - \mathbf{m\_dist}(D) + \mathbf{m\_wt}(D)\} \end{aligned}$$

where  $h_{max}$  bounds the maximum Hamming distance

CCC/ATT/AGC/CCT/ATT/GAG/ACT/GTA/CCA/GTA/AAA/TTA/AAG/CCA/GGA/ATG/GAT/GGC/CCA/AAA/GTT/AAA/CAA/TGG/CCA/TTG/ACA/GAA/GAA/AAA/ATA/AAA/GCA/TTA/GTA/GAA/ATT/TGT/ACA/GAG/ATG/GAA/AAG/GAA/GGG/\*AA/ATT/TCA/AAA/ATT/GGG/CCT/GAA/AAT/CCA/TAC/AAT/ACT/CCA/GTA/TTT/GCC/ATA/AAG/AAA/AAA/GAC/AGT/ACT/AAA/TGG/AGA/AAA/TTA/GTA/GAT/TTC/AGA/GAA/CTT/AAT/AAG/AGA/ACT/CAA/GAC/TTC/TGG/GAA/GTT/CAA/TTA/GGA/ATA/CCA/CAT/CCC/GCA/GGG/TAA/AAA/AAG/AAA/AAA/TCA/GTA/ACA/GTA/CTG/GAT/GTG/GGT/GAT/GCA/TAT/TTT/TCA/GTT/CCC/TTA/GAT/GAA/GAC/TTC/AGG/AAG/TAT/ACT/GCA/TTT/ACC/ATA/CCT/AGT/ATA/AAC/AAT/GAG/ACA/CCA/GGG/ATT/AGA/TAT/CAG/TAC/AAT/GTG/CTT/CCA/CAG/GGA/TGG/AAA/GGA/TCA/CCA/GCA/ATA/TTC/CAA/AGT/AGC/ATG/ACA/AAA/ATC/TTA/GAG/CCT/TTT/AGA/AAA/CAA/AAT/CCA/GAC/ATA/GTT/ATC/TAT/CAA/TAC/ATG/GAT/GAT/TTG/TAT/GTA/GGA/TCT/GAC/TTA/GAA/ATA/GGG/GAG/CAT/AGA/ACA/AAA/ATA/GAG/GAG/CTG/AGA/CAA/CAT/CTG/TTG/AGG/TGG/GGA/CTT/ACC/ACA/CCA/GAC/AAA/AAA/CAT/CAG/AAA/GAA/CCT/CCA/TTC/CTT/TGG/ATG/GGT/TAT/GAA/CTC/CAT/CCT/GAT/AAA/TGG/ACA/GTA/CAG/CCT/ATA/GTG/CTG/CCA/GAA/AAA/GAC/AGC/TGG/ACT/GTC/AAT/GAC/ATA/CAG/AAG/TTA/GTG/GGG/AAA/TTG/AAT/TGG/GCA/AGT/CAG/ATT/TAC/CCA/GGG/ATT/AAA/GTA/AGG/CAA/TTA/TGT/AAA/CTC/CTT/AGA/GGA/ACC/AAA/GCA/CTA/ACA/GAA/GTA/ATA/CCA/CTA/ACA/GAA/GAA/GCA/GAG/CTA/GAA/CTG/GCA/GAA/AAC/AGA/GAG/ATT/CTA/TAA/GAA/CAA/GTA/CAT/GGA/GTG/TAT/TAT/GAC/CCA/TCA/AAA/GAC/TTA/ATA/GCA/GAA/ATA/CAG/AAG/CAG/GGG/CAA/GGC/CAA/TGG/ACA/TAT/CAA/ATT/TAT/CAA/GAG/CCA/TTT/AAA/AAT/CTG/AAA/ACA/GGA/AAA/TAT/GCA/AGA/ATG/AGG/GGT/GCC/CAC/ACT/AAT/GAT/GTA/AAA/CAA/ATA/ACA/GAG/GCA/GTG/CAA/AAA/ATA/ACC/ACA/GAA/AGC/ATA/GTA/ATA/TGG/TGA/AAG/ACT/CCT/AAA/TTT/AAA/CTG/CCC/ATA/CAA/AAG/GAA/ACA/TGG/GAA/ACA/TGG/TGG/ACA/GAG/TAT/TGG/CAA/GCC/ACC/TGG/ATT/CCT/GAG/TGG/GAG/TTT/GTT/AAT/ACC/CCT/CCC/ATA/GTG/AAA/TTA/TGG/TAC/CAG/TTA/GAG/AAA/GAA/CCC

Figure 3: The genomic sequence of HIV reverse transcriptase extracted from HIV patient “AA.” Each letter (A, C, G, T) represents a nucleotide; \* represents any nucleotide. Each group of three letters represents a codon, set apart by slashes. This sequence encodes a 3D protein structure similar to that shown in Figure 1, but differing from it to some extent as specified by mutations in the sequence.

considered and  $M(S, h)$  is the set of mutants of  $S$  at Hamming distance  $h$ .  $\mathbf{m\_dist}(D)$  is the minimum Hamming distance at which a mutant occurs that resists every drug in  $D$ , and  $\mathbf{m\_wt}(D)$  is the rule weight of the least resisted drug in  $D$  by the most resistant such mutant.  $MutScore$  is comparable between drug combinations with different numbers of drugs.  $MutScore(D)$  is zero if no mutant within Hamming distance  $h_{max}$  of  $S$  resists every drug in  $D$ . Otherwise, its integer part is  $h_{max}$  minus the Hamming distance to such a mutant, and its fractional part is the maximum minimum rule weight of such mutants.

Under this model, a mutant resists a drug combination only as strongly as it resists the least-resisted drug in the combination, and a drug combination suppresses a virus population only as strongly as it suppresses the most-resistant member of the population. The effect of this is to identify nearby mutants that resist every drug in a combination, and drug combinations such that no nearby mutant resists every drug.

### Rank Alternatives

CTSHIV ranks alternative drug combinations using the current resistance weight ( $CurrWt$ ) and the nearby mutant resistances ( $MutScore$ ). The best ranked combinations of 1, 2, 3, and 4 drugs are generated independently. This is done by sorting the combinations by any monotonic function of  $CurrWt$  and  $MutScore$ . Currently we use Euclidean distance,  $\sqrt{CurrWt^2(D) + MutScore^2(D)}$ , to rank drug combination  $D$ . Values near or at zero indicate little or no resistance, and increasing positive values indicate increasing resistance. The best ranked combinations represent a satisfying compromise along both metrics simultaneously.

### Suggest Clinical Treatment Protocols

The final result of application processing is to recommend the five highest-ranked combinations of 1, 2, 3, and 4 drugs. The next highest-ranked RT-only combination is shown for comparison. Figure 4 shows 3-drug combinations recommended for an HIV patient. Figure 5 shows an example nearby resistant mutant.

It is hoped that the CTSHIV output will increase patient adherence, by clearly showing the deleterious effects of failing to take all medication. Figure 6 shows the projected consequences of non-adherence to the highest-ranked 3-drug combination of Figure 4.

### Uses of AI Technology

The key enabling AI technology is knowledge representation of the relevant scientific literature about HIV drug resistance as a set of sequence pattern rules on the HIV genome. Rule-based expert systems declaratively represent knowledge of a specialized problem and facts about a specific case, and from these draw inferences about the case. Here, the rules encode information on drug resistant mutations of HIV, the facts are the sequences of HIV genome obtained from a specific individual, and the inference to be drawn is a set of drug combinations to be recommended for the patient.

Rule forward chaining from the patient’s current HIV sequences yields currently resistant HIV mutants. Rule backward chaining through sequence space yields the nearby putatively resistant mutants. Together, they allow CTSHIV to avoid both sets of mutants.

The intelligent agent paradigm proved useful as an organizing principle. Except for the lowest level (domain-specific), Figure 2 could represent any intelligent agent connecting perception to action. Also, AI heuristic search methods are used to search sequence space.

| These protocols with 3 drugs are recommended: | CurrWt | MutScor | 0 Mut | 1 Mut | 2 Mut | 3 Mut |
|---|--------|---------|-------|-------|-------|-------|
| A5 SAQUINAVIR NELFINAVIR D4T:                 | 0.06   | 0.1     | 0.0   | 0.0   | 0.0   | 0.1   |
| B3 SAQUINAVIR DELAVIRDINE D4T:                | 0.00   | 0.2     | 0.0   | 0.0   | 0.0   | 0.2   |
| C3 SAQUINAVIR NEVIRAPINE D4T:                 | 0.00   | 0.4     | 0.0   | 0.0   | 0.0   | 0.4   |
| D4 SAQUINAVIR DELAVIRDINE AZT:                | 0.00   | 0.6     | 0.0   | 0.0   | 0.0   | 0.6   |
| E4 SAQUINAVIR NEVIRAPINE AZT:                 | 0.00   | 0.6     | 0.0   | 0.0   | 0.0   | 0.6   |
| RF3 DELAVIRDINE DDI AZT:                      | 0.08   | 1.2     | 0.0   | 0.0   | 0.2   | 0.9   |

Figure 4: Example 3-drug output from HIV patient “AA,” showing a favorable resistance profile. For the highest-ranked treatment, current resistance (CurrWt) and nearby mutation score (MutScor) are small, and only a weakly-resistant mutant appears even out to Hamming distance three (3 Mut). The letters A–F identify treatments. Treatment F is the best RT-only treatment (indicated by the prefixed letter R). Digits after the letters indicate cost codes (0 = \$0 to \$200, ..., 3 = \$600 to \$800, 4 = \$800 to \$1000, 5 = \$1000 to \$1200, ..., per month estimated average wholesale cost).

|  | CurrWt | MutScor | 0 Mut | 1 Mut | 2 Mut | 3 Mut |
|--|--------|---------|-------|-------|-------|-------|
| A5 D4T NELFINAVIR SAQUINAVIR:  | 0.06   | 0.1     | 0.0   | 0.0   | 0.0   | 0.1   |
| Current: (NELFINAVIR) RT 151:CAG->ATG by R11 (D4T) PRO 90:TTG->ATG by R28 (SAQUINAVIR) |        |         |       |       |       |       |

Figure 5: Example output for HIV patient “AA” showing one example of the closest mutants inferred to most resist every drug in the top-ranked 3 drug combination of Figure 4. Three letters must change simultaneously. Currently Nelfinavir is resisted; changing two letters at RT 151 resists D4T and changing one at PRO 90 resists Saquinavir.

|                               | CurrWt | MutScor | 0 Mut | 1 Mut | 2 Mut | 3 Mut |
|-------------------------------|--------|---------|-------|-------|-------|-------|
| A5 SAQUINAVIR NELFINAVIR D4T: | 0.06   | 0.1     | 0.0   | 0.0   | 0.0   | 0.1   |
| If stop NELFINAVIR:           |        | 0.6     | 0.0   | 0.0   | 0.0   | 0.6   |
| If stop SAQUINAVIR:           |        | 1.1     | 0.0   | 0.0   | 0.1   | 1.0   |
| If stop D4T:                  |        | 2.1     | 0.0   | 0.1   | 0.6   | 1.1   |

Figure 6: Example output for HIV patient “AA” showing the projected result from stopping any single drug in the top-ranked 3 drug combination of Figure 4 (Saquinavir, Nelfinavir, D4T). Mutants are closer or worse or both. Stopping Nelfinavir is bad, stopping Saquinavir is worse, and stopping D4T is worst of all.

## Application Use and Payoff

The first HIV patient data was run through the CTS-HIV system in June, 1996. In February, 1997, the application began its first round of human clinical trials on 15 HIV patients at the University of California, Irvine, and at the Center for Special Immunology as a satellite site. Informed consent was obtained using a form approved by the UC Irvine Institutional Review Board. All patients had detectable viral load at baseline (mean  $\log_{10}$  load of  $4.67 \pm 2.16$ ), weakened immune system (CD4 counts  $< 500$  cells/mm<sup>3</sup>), and failure of at least one previous antiviral treatment regimen.

While these trials are still ongoing, initial results are encouraging (See *et al.* 1998). At three months, 10 of 15 patients who had failed at least one prior treatment regimen had an undetectable viral load (67% success). This compares to about 20% in everyday practice in the same patient population. Currently, at the midway point of the study, 12 of 15 patients have undetectable viral loads at six or nine months (80% success). The other three subjects all have at least a 2-log reduction in viral load compared to baseline. Note that in usual clinical trials, the percentage of viral-load undetectable

patients diminishes over time. We expect and are seeing *improvement* over time based upon CTSHIV suggested treatment regimens. Further detailed results will be presented in the domain literature.

Currently a total of 58 HIV patients have been run through the CTSHIV system. A new round of phase II clinical trials of CTSHIV involving 30 HIV patients who previously failed multiple antiretroviral regimens has been initiated at the University of California, Irvine, and Stanford University, Santa Clara campus. The control group will use plasma HIV RNA and CD4 cell counts as in conventional therapy, while the test group will be identical but also provide CTSHIV recommendations to the primary care physician. Collaborations with several other groups involved in the treatment of HIV patients have begun and are expanding. An Affymetrix gene chip machine has been purchased and sequencing throughput will increase dramatically when it comes online. Because of the early encouraging results of the clinical trials, wide-spread recognition of the drug resistance problem, and the high rate of HIV infection in the general population, we expect use of the application to increase sharply in the near future.

## Application Development, Deployment

Three domain experts (Darryl See, Douglas Richman, Edison Schroeder) began extracting rules from the scientific literature in September, 1995. The first rule set was completed in May, 1996.

The first rule-based system prototype was developed to identify current resistance already present in the patient's HIV infection (Pazzani *et al.* 1997b). It was coded in FOCL-1-2-3 (Pazzani & Kibler 1992), a LISP based expert shell. It was begun in March, 1996, and completed in June, 1996. It was re-coded in JAVA between April and June of 1997 (Pazzani *et al.* 1997a).

The ability to use the rules to search mutation sequence space for nearby drug resistant mutants was unanticipated when the original knowledge was encoded and the first prototype developed, and so demonstrates the robustness and extensibility of knowledge-based systems. A LISP based mutation space search engine was begun in November, 1996, and completed in May, 1997. The two subsystems were integrated and re-coded in LISP between October and December, 1997.

The application is deployed primarily by the email exchange of input clinical data and output recommended treatments. We have developed an automatic email server, as well as a WWW-based graphical interface to the email server. The server extracts patient data from the body of an email message, automatically enqueues the application to process it, and emails the results back to the sender.

Deployment has been smooth largely because the application end-users so far have been enthusiastic domain experts who are currently treating HIV patients. For cases where a treatment regimen has failed due to the development of drug resistance, the application enables them to base their next choice of treatment regimen on scientific principles and experimental data. This replaces the blind intuition and guess-work that formerly guided treatment switches after treatment failure. They are glad to see their patients improve, anxious to see the application succeed, and tolerant of the few glitches.

## Maintenance

It is doubtful that the knowledge base will be complete until HIV is eradicated. Maintenance of CTSHIV is equivalent to adding new rules from the scientific AIDS literature. The rules are revised by three domain experts every three months by extracting new rules that have appeared in the literature in the interim. Relevant articles are retrieved by keyword-based literature search, old rules revised as needed, and new rules composed manually.

In the future we anticipate that the challenge of extending the knowledge base will provide fruitful opportunities for intelligent applications. An intelligent information retrieval system could monitor the literature, retrieve papers that mention HIV drug resistant mutations, extract candidate rules, and automatically enqueue review by domain experts. Other AI approaches

could suggest when to test a patient strain for possible resistance to a specific drug. Predicting when a putative mutant is unviable, and coping with resistance that occurs outside the rule set, are further challenges for intelligent systems. Machine learning and data mining techniques could learn new rules, infer trends and recognize regularities in resistance patterns.

## Summary

We have described an AI application (CTSHIV) that connects the scientific literature describing specific HIV drug resistances directly to the HIV virus strain infecting a specific HIV patient. The application identifies current and nearby drug resistant mutant strains, ranks the current FDA-approved treatment regimens according to their estimated ability to avoid the resistant mutants, and recommends a Customized Treatment Strategy for the individual patient involved. Thus the significance of the application is (1) a method for addressing HIV drug resistance in the clinic, especially treatment switches after treatment failure, based on scientific principles and experimental data, (2) a decreased tendency to select for drug resistance in the global HIV gene pool, and (3) a possible model for the use of knowledge-based systems in other drug resistant viruses.

This paper also illustrated the robustness and extensibility of knowledge-based systems. It showed how knowledge originally encoded to perform one knowledge-based task easily may be re-directed to perform another, even one not anticipated when the original knowledge was encoded. This result supports knowledge-base efforts to encode knowledge in socially important areas.

## Acknowledgments

Douglas Richman and Edison Schroeder helped extract the original knowledge base from the AIDS scientific literature, and help revise the knowledge base every 3 months. Doug Cable and Winnie Huang co-supervised the first CTSHIV clinical trials. Data entry and analysis was done by Richard Haubrich and Allen McCutchan at the University of California at San Diego. Carol Kemper is co-supervising the next round of clinical trials. Tom Gingeras contributed useful practical insights and advice. We gratefully acknowledge the technical assistance of Tonya Clark and Marikel Chatard in cloning, PCR, and RNA/DNA extraction.

Funding has been provided by Roche Molecular Systems, the California University-wide AIDS Research Program through the California Collaborative Treatment Group (CCTG), and the National Science Foundation under grant IRI-9624739.

CTSHIV is available from University/Industry Research and Technology, 380 University Tower, University of California, Irvine, 92717 USA.

## References

- Carpenter, C.; Fischl, M.; Hammer, S.; Hirsch, M.; Jacobsen, D.; Katzenstein, D.; Montaner, J.; Richman, D.; Saag, M.; Schooley, R.; Thompson, M.; et al. 1996. Antiretroviral therapy for HIV infection in 1996. *J. American Medical Assoc.* 276:146–154.
- CDC. 1997. *HIV/AIDS Surveillance Report*. Centers for Disease Control and Prevention, Atlanta, Georgia, USA, 9(1) edition.
- Condra, J.; Schlieff, W.; Blahy, O.; Gabryelski, L.; Graham, D.; Quintero, J.; Rhodes, A.; Robbins, H.; Roth, E.; Shivaprakash, M.; Titus, D.; et al. 1995. *In vivo* emergence of HIV-1 variants resistant to multiple protease inhibitors. *Nature* 374:569–571.
- Diamond, L.; Nguyen, D.; Jouault, H.; Imbert, M.; and Sultan, C. 1994. An expert system for the interpretation of flow cytometric immunophenotyping data. *J. of Clinical Computing* 22:50–58.
- Finzi, D.; Hermankova, M.; Pierson, T.; Carruth, L.; Buck, C.; Chaisson, R.; Quinn, T.; Chadwick, K.; Margolick, J.; Brookmeyer, R.; Gallant, J.; et al. 1997. Identification of a reservoir for HIV-1 in patients on highly active antiretroviral therapy. *Science* 278:1295–1300.
- Gu, Z.; Gao, Q.; Fang, H.; Parniak, M.; Brenner, B.; and Wainberg, M. 1994. Identification on novel mutations that confer drug resistance in the human immunodeficiency virus polymerase gene. *Leukemia* 8S1:5166–5169.
- Hsiou, Y.; Ding, J.; Das, K.; Hughes, S.; and Arnold, E. to be published. Structure of unliganded human immunodeficiency virus type 1 reverse transcriptase at 2.7 Angstroms resolution.
- Iversen, A.; Shafer, R.; Wearnly, K.; Winters, M.; Mullins, J.; Chesebro, B.; and Morgan, T. 1996. Multidrug-resistant human immunodeficiency virus type 1 strains resulting from combination antiretroviral therapy. *J. Virology* 70:1086–1090.
- Jacobsen, H.; Hanggi, M.; Ott, M.; Duncan, I.; Owen, S.; Andreoni, M.; Vella, S.; and Mous, J. 1996. *In vivo* resistance to a human immunodeficiency type-1 protease inhibitor. *J. Infect. Diseases* 173:1379–1387.
- Lange, J. 1995. Triple combinations: present and future. *J. of AIDS and Human Retrovirology* 10 Suppl 1:S77–82.
- Musen, M.; Wieckert, K.; Miller, E.; Campbell, K.; and Fagan, L. 1995. Development of a controlled medical terminology: knowledge acquisition and knowledge representation. *Methods of Information in Medicine* 34:85–95.
- Musen, M.; Tu, S.; Das, A.; and Shahar, Y. 1996. EON: a component-based approach to automation of protocol-directed therapy. *J. Amer. Medical Informatics Assoc.* 3:367–388.
- O'Brian, W.; Hartigan, P.; Martin, D.; Esinhart, J.; Hill, A.; Benoit, S.; Rubin, M.; Simberkoff, M.; and Hamilton, J. 1996. Changes in plasma HIV-1 and CD4+ lymphocyte counts and the risk of progression to AIDS. *New England J. Medicine* 334:426–431.
- Ohno-Machado, L.; Parra, E.; Henry, S.; Tu, S.; and Musen, M. 1993. AIDS 2: A decision-support tool for decreasing physician's uncertainty regarding patient eligibility for HIV treatment protocols. In *Proc. of the 17th Annual Symp. on Computer Applications in Medical Care*, 429–433.
- Pazzani, M., and Kibler, D. 1992. The utility of prior knowledge in inductive learning. *Machine Learning* 9:54–97. FOCL-1-2-3 is available from <http://www.ics.uci.edu/~mlearn/FOCL.html>.
- Pazzani, M.; Iyer, R.; See, D.; Shroeder, E.; and Tilles, J. 1997a. CTSHIV: A knowledge-based system in the management of HIV-infected patients. In *Proc. of the Intl. Conf. on Intelligent Information Systems*.
- Pazzani, M.; See, D.; Shroeder, E.; and Tilles, J. 1997b. Application of an expert system in the management of HIV-infected patients. *J. of AIDS and Human Retrovirology* 15:356–362.
- Ruggiero, C.; Giacomini, M.; Varnier, O. E.; and Gaglio, S. 1994. A qualitative process theory based model of the HIV-1 virus-cell interaction. *Computer Methods and Programs in Biomedicine* 43:255–259.
- Safran, C.; Rind, D.; Sands, D.; Davis, R.; Wald, J.; and Slack, W. 1996. Development of a knowledge-based electronic patient record. *M.D. Computing* 13:46–54.
- See, D.; Cimoch, P.; Lathrop, R.; Pazzani, M.; and Reiter, W. 1998. Successful use of an expert program in managing antiretroviral use in patients that have failed haart therapy. In *Proc. of the 9th Intl. AIDS Conf.*, to appear.
- Sieburg, H. 1994. Methods in the Virtual Wetlab I: rule-based reasoning driven by nearest-neighbor lattice dynamics. *AI in Medicine* 6:301–319.
- Sobesky, M.; Michelet, C.; Thomas, R.; and LeBeux, P. 1994. Decision making system. *J. Clinical Computing* 22:20–26.
- Sonnenberg, F.; Hagerty, C.; and Kulikowski, C. 1994. An architecture for knowledge-based construction of decision models. *Medical Decision Making* 14:27–39.
- Tu, S.; Eriksson, H.; Gennari, J.; Shahar, Y.; and Musen, M. 1995. Ontology-based configuration of problem-solving methods and generation of knowledge-acquisition tools. *AI in Medicine* 7:257–289.
- Wong, J.; Hezareh, M.; Günthard, H.; Havlir, D.; Ignacio, C.; Spina, C.; and Richman, D. 1997. Recovery of replication-competent HIV despite prolonged suppression of plasma viremia. *Science* 278:1291–1295.
- Xu, L. 1996. An integrated rule- and case-based approach to AIDS initial assessment. *Intl. J. of Bio-Medical Computing* 40:197–207.