

Game Theoretic Antibody Design

Swetasudha Panda
Electrical Engineering and Computer
Science
Vanderbilt University
Nashville, Tennessee

Alexander M. Sevy
Center for Structural Biology
Vanderbilt University
Nashville, Tennessee

James E. Crowe Jr
Vanderbilt Vaccine Center
Vanderbilt University Medical Center
Nashville, Tennessee

Jens Meiler
Center for Structural Biology,
Department of Chemistry
Vanderbilt University
Nashville, Tennessee

Yevgeniy Vorobeychik
Electrical Engineering and Computer
Science
Vanderbilt University
Nashville, Tennessee

ABSTRACT

Vaccines have had a marked impact on public health, in some cases eradicating previously common debilitating or deadly diseases, such as polio. However, designing vaccines is a laborious processes, involving extensive and expensive experimentation, and a great deal of trial and error. Computational vaccine design promises to be game changing in this regard. Such design commonly involves two steps: computational design of antibodies, followed by computational design of a vaccine molecule which would promote generation of such antibodies. We focus on the first step of this pipeline: antibody design.

An important challenge in modern antibody design is the possibility of rapid viral mutations which escape antibody binding, as is the case with HIV and influenza. Indeed, in both these cases, evolution of the viral antigen has thus far foiled attempts at designing an effective long-term vaccine. A common way aimed at capturing viral evolution is to use a *fixed* panel of known viral variants, with the goal of designing a *broadly binding* antibody (i.e., one which binds most, or all of these). We propose a novel game theoretic approach to this problem, which allows us to capture not merely a fixed panel of viral variants, but also a combinatorial space of mutations from these. Our approach combines learning a linear approximation of binding stability energy of the antibody-virus complex with bi-level integer linear programming, which we transform into a single-level mixed-integer linear program. Through a series of simulation experiments we demonstrate the efficacy of our proposed approach.

KEYWORDS

Stackelberg games; Combinatorial optimization; Machine learning; Applications in computational biology

1 INTRODUCTION

Infectious diseases pose a major threat to public health. In 2016, about 36.7 million people were living with HIV, and it resulted in

1 million deaths [14]. From the time AIDS was identified, it has caused an estimated 35 million deaths worldwide [20]. A recent Ebola outbreak in Africa killed thousands [3], and annual influenza outbreaks affect millions, with hundreds of thousands hospitalized, and thousands dying from the influenza or its side-effects [4].

Vaccination therapies are among the most important methods for combating infectious diseases. Vaccines are external substances that stimulate the immune system to produce antibodies that bind to the vaccine substance. As antibodies develop in response to a vaccine against a particular pathogen, they remain in the individual's bloodstream and rapidly neutralize and clear the pathogen if the individual is ever infected, thereby preventing illness. Traditional vaccine design involves laborious and costly lab work aimed at finding just the right substance which would successfully and reliably elicit antibodies binding the target pathogen. Recently, a promising approach has been taking shape in which vaccines are designed *computationally*, making use of modern computational protein modeling tools, such as ROSETTA [1]. One of the common approaches involves two steps: first, finding an antibody with desired neutralization characteristics, and second, finding a vaccine which binds tightly to the desired antibody, thereby eliciting the associated target immune response. We focus on the first step of computational antibody design.

The central goal in computational antibody design is to find an antibody protein sequence which neutralizes the target pathogen. In order for the antibody to neutralize a pathogen, it needs to *bind* to it; the specific position at which the two typically bind is called the *binding site*. When two proteins (such as an antibody and viral proteins) bind, they form a *complex*, which is a configuration minimizing the total energy of the two molecules. Binding is typically highly specific: a small change in the sequence can destabilize binding.

However, binding a single fixed antigen (portion of the pathogen which typically interacts with the antibody) is often insufficient: for example, viruses such as HIV and flu have many strains, and an antibody which neutralizes one will often fail to neutralize another. An area of active research in antibody design (computational and otherwise), therefore, is to develop and characterize *broadly binding antibodies*, that is, antibodies which effectively bind to (and, ideally, neutralize) many variants of the pathogen [10]. Nevertheless,

as a pathogen evolves, it may well still escape neutralization; for example, HIV has an extremely high mutation rate [8].

We propose a radically different approach for computational antibody design in the context of rapidly mutating viruses: using a game theoretic (Stackelberg game) model for the interaction between the antibody and the virus. In this game, the antibody designer chooses an antibody sequence, while the virus aims to maximally destabilize binding to the resulting antibody, subject to a constraint on the number of mutations (this constraint captures the fact that such a mutation has to be sufficiently likely). This game can be formulated as a bi-level optimization problem; unfortunately, such a formulation is quite intractable. We address tractability in three steps: first, we learn a linear approximation of the antibody-virus binding score as a function of its sequence (including all pairwise interactions at the binding site); second, we formulate the optimal virus escape problem as an integer linear program; and third, after relaxing the integrality constraint in the virus escape program and taking its dual, we formulate the antibody design bi-level problem as a mixed-integer linear program. Our experimental results demonstrate that our approach is extremely effective against two recent prior approaches for HIV antibody design.

2 RELATED WORK

Conceptually, our work follows the research on game theoretic antibody design in [15, 16], which extends the insights of Stackelberg security games [11] to the vaccine design domain. This previous game-theoretic approach has a fundamental limitation as it relies on local search approaches. Typically, it is extremely challenging to compute an optimal solution to bi-level problems with integer variables. The most important contribution in this paper is that the compact formulation allows us to compute the optimal global solution. Additionally, our model incorporates both binding and stability energies that are critically important factors in protein sequence design. Since antibodies are protein sequences, our work relates significantly to computational protein design. Recent advances involve multi-specificity design to achieve protein design with respect to more than one targets [18]. The most relevant prior work is Breadth Optimization in Antibody Design (BROAD) that incorporates machine learning and sequence optimization for efficient sampling in the sequence space [19]. However, while BROAD maximizes breadth over an existing virus panel, our approach of game-theoretic antibody design goes significantly further as the designed antibody continues to bind against virus escape mutations. There have been numerous efforts in learning protein structure, function and interactions from sequence data, of which Kamisetty et al. [12] is the most relevant to our effort. More remotely related work include game theoretic models of vaccination decisions [2, 5, 13]. However, these model human decisions about being vaccinated, whereas our model involves molecular-level interactions between immunity and pathogen.

3 A GAME THEORETIC MODEL OF ANTIBODY DESIGN

We define an antibody or virus primary sequence as a sequence (vector) of amino acids as in previous work [16]. Let c denote the native virus (the initial virus strain before mutations) and (a, v) be

arbitrary antibody and virus sequences respectively. Let $\mathcal{B}(a, v)$ and $\mathcal{S}(a, v)$ denote the binding energy and the thermodynamic stability scores of the antibody-virus complex. A combination of these is used as the overall *energy score* (often known as the z-score) of the complex, which is what we actually work with, and denote by $\mathcal{Z}(a, v)$. Also, lower (more negative) scores indicate stronger binding and stability of the antibody-virus complex.

The virus sequence attempts to escape binding to the antibody by making a series of mutations. We can represent the number of mutations in a mutated virus sequence v from the native c as $\|v - c\|_0$, where the l_0 norm computes the number of sequence positions in v that are different from c . Given an antibody a , we model the virus as making up to α mutations with the goal of maximizing its binding energy score so as to destabilize binding. This model is motivated by natural selection: viral proteins which tightly bind to an antibody will be cleared by the immune system, leaving those which do not, and the remaining viral variants, mutating from a native sequence, will thereby increase in relative prevalence.

In general, there are many potential virus variants that can infect an individual. To capture this, we consider T virus sequences of different types t in a virus panel, each starting from a native sequence c^t and making mutations to escape binding to a .

We can formally represent the optimization problem being solved by a collection of viruses as follows:

$$\begin{aligned} & \underset{v^t \in \mathcal{V}}{\text{maximize}} && \sum_{t=1}^T \mathcal{Z}(a, v^t) \\ & \text{subject to} && \|v^t - c^t\|_0 = \alpha, \forall t. \end{aligned} \quad (1)$$

where \mathcal{V} is the space of virus sequences under consideration. The optimization problem (1) can be viewed as the combined *best response* of the virus panel to a fixed antibody a .

The space of feasible virus sequences \mathcal{V} can be all possible combinations of amino acids in corresponding positions. However, in practice many such combinations are infeasible in nature, for example, because some mutations in specific positions destabilize the viral protein, or affect function. These considerations are too complex to capture cleanly. As a proxy, we constrain feasible mutations in each position to those which have been observed in nature (in that position) sufficiently often. More precisely, we only consider a mutation in a position i to an amino acid j if $p_{ij} \geq \theta$, where p_{ij} is the empirical frequency of the associated position-specific mutation, and θ an exogenously specified threshold ($\theta = 0$ is a natural choice, and one we use in the experiments; at this threshold, we only disallow mutations that have *never* been observed in nature).

In addition to only allowing mutations which are not too rare in nature, we impose another natural restriction on \mathcal{V} . Specifically, first-order effects in regard to its antibody binding properties are determined by the sequence that is a part of the native virus *binding site* (i.e., positions on the native virus sequence which are in contact with the native antibody in the original binding complex). Therefore, we only consider the problem of virus escape in terms of binding site mutations. This also allows us to significantly reduce the dimensionality of the problem in practice.

Now we consider the problem of designing an antibody, a , that is robust to virus escape, as we have now formally defined using the optimization problem (1). The antibody designer's decision

problem is then to choose an antibody which minimizes the energy scores (strengthens binding and stability) *with respect to the virus panel* $\{1, \dots, T\}$, accounting for potential mutations of each virus in response. This gives rise to the following bi-level optimization problem for antibody design:

$$\begin{aligned} \min_{\mathbf{a} \in \mathcal{A}} \max_{\mathbf{v}^t \in \mathcal{V}} \sum_{t=1}^T \mathcal{Z}(\mathbf{a}, \mathbf{v}^t) \\ \text{subject to } \|\mathbf{v}^t - \mathbf{c}^t\|_0 = \alpha, \forall t \end{aligned} \quad (2)$$

where \mathcal{A} is the antibody design space, which we restrict to the native binding site for the same reasons as for the virus. Observe that the antibody-virus interaction in our model can be viewed as a Stackelberg game in which the designer (antibody) is the leader, and each virus is the follower, who chooses an alternative virus sequence in response to the antibody chosen by the designer. Moreover, this game is zero-sum: the designer minimizes the energy score, a quantity which is maximized by each virus t .

4 SOLUTION APPROACH

4.1 A Bi-Linear Representation of Energy Scores

The optimization problem (2) is intractable in general, even when simulated using the ROSETTA software. In particular, computing such a function using ROSETTA even for a given pair of sequences requires many runs of stochastic local search, and takes on the order of minutes or hours. We make progress by approximating the complex black-box ROSETTA energy function $\mathcal{Z}(\mathbf{a}, \mathbf{v})$ by a bi-linear function of the antibody and virus sequences, similar to the approach proposed by Kamisetty et al. [12]. The model is based on an assumption that the binding and stability of an antibody-virus complex is primarily determined by two factors: a) the individual amino acids in each binding position of the antibody and the virus respectively, and b) the effects of the pairwise amino acid interactions between the antibody and the virus. We now describe this model in detail.

We represent an antibody sequence \mathbf{a} as a binary position by amino-acid matrix, with $a_{ij} = 1$ iff amino acid j appears in position i , and $a_{ij} = 0$ otherwise. Thus, $\sum_j a_{ij} = 1$, since exactly 1 amino acid can be in a given position. Similarly, the virus protein sequence is represented as a binary matrix v_{ij} which is 1 iff amino acid j is in position i . Let N_a and N_v denote the number of binding positions on the antibody and the virus respectively, and let $M = 20$ denote the number of amino acids.

Amino acid contributions to the energy score can be modeled as a bipartite graph in which nodes represent the amino acids and the edges represent the pairwise amino acid interactions. Each antibody position node i has an associated weight vector $\mathbf{x}_i \in \mathbb{R}^M$. Similarly, each virus position node j has an associated weight vector $\mathbf{y}_j \in \mathbb{R}^M$.

The edge (i, j) between antibody position node i and virus position node j has an associated weight matrix $Q_{ij} \in \mathbb{R}^{M \times M}$ to represent the position specific contribution to the energy score for each amino acid pair. Consequently, given \mathbf{a} and \mathbf{v} , the energy score varies as the sum of individual amino acids and pairwise interaction effects. Given this setting, the z-score for a given pair \mathbf{a} and \mathbf{v} is

defined as:

$$\begin{aligned} \mathcal{Z}(\mathbf{a}, \mathbf{v}) = & \sum_{i=1}^{N_a} \sum_{j=1}^M x_{ij} a_{ij} + \sum_{i=1}^{N_v} \sum_{j=1}^M y_{ij} v_{ij} \\ & + \sum_{k=1}^{N_a} \sum_{l=1}^{N_v} \sum_{u=1}^M \sum_{m=1}^M a_{ku} q_{kl}^{um} v_{lm} + I \end{aligned} \quad (3)$$

where I is the intercept term and q_{kl}^{um} represents $Q_{kl}(u, m)$.

Our bi-linear model thus has four sets of parameters: \mathbf{x}_i , \mathbf{y}_j , and Q_{ij} for all pairs of antibody and virus positions, i and j , respectively, and the intercept I . We learn these parameters by generating a dataset of ROSETTA energy function values for a number of pairs of antibody and virus sequences (as detailed in the experiments).

Armed with the bi-linear model described in this section, we can convert the hard bilevel optimization problem into a significantly more tractable mixed-integer linear program through a combination of convex relaxation and duality, as we describe next.

4.2 Integer Linear Program for Virus Escape

Our first step is to formulate the virus optimal escape problem as an integer linear program.

We start by observing that the number of mutations α can be computed using a dot product with the sequence representation described above. Specifically, $\mathbf{v}^t \cdot \mathbf{v}^t = N_v$ and $\mathbf{v}^t \cdot \mathbf{c}^t = N_v - \alpha$. Moreover, $\mathcal{Z}(\mathbf{a}, \mathbf{v})$ is now a linear function with the above sequence representation. These observations allow us to formulate the virus escape optimization in Equation 1 as an integer linear program (ILP). Since in this problem the antibody \mathbf{a} is fixed, we can group the model in Equation 3 in terms of the variables \mathbf{v} as $\sum_{i=1}^{N_a} \sum_{j=1}^M x_{ij} a_{ij} + \sum_{i=1}^{N_v} \sum_{j=1}^M (y_{ij} + \sum_{k=1}^{N_a} \sum_{u=1}^M a_{ku} q_{kl}^{ki}) v_{ij} + I$. Thus, the virus escape ILP for a particular native virus indexed by t (from a collection of T of these) can be formulated as follows:

$$\begin{aligned} \text{maximize}_{\mathbf{v}^t \in \mathcal{V}} \sum_{i=1}^T \sum_{j=1}^{N_v} \sum_{k=1}^M (y_{ij} + \sum_{k=1}^{N_a} \sum_{u=1}^M a_{ku} q_{kl}^{ki}) v_{ij}^t \end{aligned} \quad (4a)$$

$$\begin{aligned} & + T \sum_{i=1}^{N_a} \sum_{j=1}^M x_{ij} a_{ij} \\ \text{subject to } \sum_{j=1}^M v_{ij}^t &= 1, \forall i, t \end{aligned} \quad (4b)$$

$$N_v - \sum_{i=1}^{N_v} \sum_{j=1}^M v_{ij}^t c_{ij}^t = \alpha, \forall t \quad (4c)$$

$$v_{ij}^t \leq L(p_{ij} - \theta), \forall i, j, t \quad (4d)$$

$$v_{ij}^t \in \{0, 1\}, \forall i, j, t$$

where constraint 4b enforces that the binary variables v_{ij}^t at each antibody binding position should sum to 1, i.e., each position admits one amino acid. The term $\sum_{i=1}^{N_v} \sum_{j=1}^M v_{ij}^t c_{ij}^t$ in constraint 4c computes the dot product $\mathbf{v}^t \cdot \mathbf{c}^t$. The constraint 4d encodes the constraint that we only allow mutations at positions to amino acids which

have been observed at a frequency $p_{ij} \geq \theta$ as a linear constraint; here, L is a large number.

4.3 Mixed Integer Linear Program for Antibody Design

While we can represent the optimization problem faced by the virus *given a fixed antibody* using a linear integer program, our underlying problem of antibody design is still a bi-level problem. Such bi-level problems (with integer variables, as in our case) are, in general, extremely challenging to solve.

At the high level, we propose to leverage the linear structure of the problem to solve it. First, we relax the integrality constraint of the inner (virus escape) problem. This yields a linear program, the dual of which we embed into the outer integer linear program. By relaxation, combined with strong duality of linear programming, the resulting mixed-integer linear program minimizes an *upper bound* on the z-score objective with respect to optimal virus escape.

We start with the ILP 4 computing the optimal virus escape, and relax the integrality constraint; that is, we relax the binary v_{ij}^t variables to be continuous and add the constraints $0 \leq v_{ij}^t \leq 1$. Next, we show that the resulting relaxed LP has integral optimal solutions.

The standard form LP $\{\max w^T s : As = b, s \geq 0\}$ with integral right-hand side vector b has an integral optimal solution if its constraint matrix A is *totally unimodular*, a notion we now define.

Definition 4.1 (Total Unimodularity). A matrix A is totally unimodular (TUM) if the determinant of each square submatrix of A is in $\{0, 1, -1\}$. In particular, each entry of A is in $\{0, 1, -1\}$.

THEOREM 4.2 (SUFFICIENT CONDITION). A matrix A is TUM if it only has at most two non-zero entries 1 or -1 in every column, and for all columns with two non-zero coefficients, the column sum is 0.

We next use this sufficient condition for TUM to prove that our LP relaxation yields optimal solutions to the original ILP. This result allows us to work with the relaxed LP for the virus escape problem.

PROPOSITION 4.3. The LP relaxation of the virus escape ILP 4 has integer optimal solutions.

PROOF. We first prove that the constraint matrix in the LP relaxation is TUM. Consider the LP relaxation with the constraints 4b and 4c and the non-negative variables. The additional constraints $0 \leq v_{ij}^t \leq 1$ are already enforced using 4b and the fact that all variables are non-negative. The corresponding constraint matrix has at most two non-zero elements in any given column corresponding to the variables v_{ij}^t . The first non-zero element +1 from the relevant constraint 4b and the second non-zero element -1 from 4c. Therefore, using theorem 4.2, the constraint matrix is TUM. Since the right hand side vector has integer elements, this LP relaxation has optimal integer solutions. This conclusion continues to hold after adding the constraints 4d since these only additionally restrict the variables to be zero under specific conditions. \square

We observe that the primal relaxed LP is feasible and bounded, and, therefore, the dual is also feasible and bounded, and (by strong duality) has the same solution as the primal. Let the associated (non-negative) dual variables be denoted by ψ_{ij}^t for each of the

constraints $v_{ij}^t \leq 1$, and let ϕ_i^t (unrestricted), π^t (unrestricted) and ρ_{ij}^t (non-negative) denote the dual variables corresponding to constraints 4b, 4c, and 4d. Note that all dual variables are continuous. The dual LP is the given by the following (a is fixed here as in the primal LP):

$$\text{minimize}_{\phi, \psi, \rho, \pi} \sum_{t=1}^T \left[\sum_{i=1}^{N_v} \phi_i^t - (N_v - \alpha)\pi^t + \sum_{i=1}^{N_a} \sum_{j=1}^M L(p_{ij} - \theta)\rho_{ij}^t \right] \quad (5a)$$

$$+ \sum_{i=1}^{N_v} \sum_{j=1}^M \psi_{ij}^t \Big] + T \sum_{i=1}^{N_a} \sum_{j=1}^M x_{ij} a_{ij}$$

$$\text{subject to } \phi_i^t - \pi^t c_{ij}^t + \rho_{ij}^t + \psi_{ij}^t \quad (5b)$$

$$\geq \left(y_{ij} + \sum_{k=1}^N \sum_{u=1}^M a_{ku} q_{uj}^{ki} \right), \forall i, j, t \quad (5c)$$

$$\psi, \rho \geq 0, \pi, \phi \text{ unrestricted variables}$$

Next, we integrate this dual LP into the antibody optimization problem in Equation 2 to formulate the following mixed integer linear program (MILP):

$$\text{minimize}_{a \in \mathcal{A}, \phi, \psi, \rho, \pi} \sum_{t=1}^T \left[\sum_{i=1}^{N_v} \phi_i^t - (N_v - \alpha)\pi^t + \sum_{i=1}^{N_a} \sum_{j=1}^M L(p_{ij} - \theta)\rho_{ij}^t \right] \quad (6a)$$

$$+ \sum_{i=1}^{N_v} \sum_{j=1}^M \psi_{ij}^t \Big] + T \sum_{i=1}^{N_a} \sum_{j=1}^M x_{ij} a_{ij}$$

$$\text{subject to } \phi_i^t - \pi^t c_{ij}^t + \rho_{ij}^t + \psi_{ij}^t \quad (6b)$$

$$- \sum_{k=1}^N \sum_{u=1}^M a_{ku} q_{uj}^{ki} \geq y_{ij}, \forall i, j, t \quad (6c)$$

$$\sum_{u=1}^M a_{ku} = 1, \forall u \quad (6d)$$

$$a_{ij}^t \in \{0, 1\}, \forall i, j, t$$

$$\psi, \rho \geq 0, \pi, \phi \text{ unrestricted variables}$$

The variables now include the binary antibody variables a_{ij} , and the constraints ensure that these sum to 1 at each antibody binding position, i.e., each position admits one amino acid. An important observation we can make is that while originally we had bi-linear terms involving antibody and virus decision variables, these are decoupled after taking the dual, resulting in solely linear terms.

5 EXPERIMENTS

We denote our proposed antibody design approach as STRONG: STackelberg game theoretic model for RObust aNtibody desiGn and compare against the two prior approaches, a) BROAD [19] and b) the game theoretic approach proposed in [16] (henceforth denoted as AAMAS2015).

The data comprises the anti-HIV antibody VRC23 [9] (the native antibody) against a set of 180 diverse HIV gp120 virus sequences (derived from Chuang et al. [6]). To generate sufficient training data that consists of antibody and virus sequence pairs and the associated scores, we make random substitutions in the binding sites of VRC23 and the set of 180 virus sequences ($N_a = 27$ and

$N_v = 32$). Each antibody/virus variant has five randomly selected amino acid substitutions. All such antibody-virus pairs are subjected to an energy minimization via the ROSETTA relax protocol (iterative rounds of side chain repacking and backbone minimization [7], talaris2013 score function). We generate 50 models of each antibody-virus pair resulting from the above energy minimization protocol and choose the lowest scoring model in each case. This allows us to construct the dataset for our experiments with a total of 7360 such random antibody-virus combinations (including VRC23 and the 180 virus sequences), and the associated scores. We compute mutation frequencies (p_{ij} in our terminology) from an exhaustive database of over 66,000 HIV-1 sequences (from the Los Alamos HIV sequence database <http://www.hiv.lanl.gov/>). We set the threshold θ to be 0, excluding only mutations which are *never* observed in nature. Finally, we evaluate the robust antibody sequences generated using our proposed approach using simulation experiments (and ROSETTA structure modeling for a subset of experiments because of the high computational cost).

5.1 Bi-linear Z-score Model

The feature vector \mathbf{f} consists of $N_a \times M$ binary antibody features, $N_v \times M$ binary virus features and $N_a \times N_v \times M \times M$ binary pairwise interaction features corresponding to \mathbf{x} , \mathbf{y} and \mathbf{Q} respectively. We use sparse matrices to represent this feature space and use the Lasso implementation in scikit-learn [17] with l_1 (sparse) regularization. To measure the accuracy of predictions, we compute the correlation coefficient between the ROSETTA computed z-scores and the scores predicted by regression. We perform a 10-fold cross validation experiment with 80% of the data for training and 20% for testing. Based on this parameter tuning, we choose regularization parameter $\lambda = 0.01$ with an average correlation of 0.85 between the predicted and the ROSETTA computed z-scores.

5.2 Comparison against BROAD

BROAD [19] is a state of the art algorithm for antibody design against a *fixed* panel of HIV virus variants that involves generating a large training set of binding and stability scores using ROSETTA, fitting linear models to predict binding and stability, and solving an ILP to compute an optimal broadly binding antibody sequence.

We perform the comparison following the experimental workflow in BROAD. We construct 50 random subsamples of the full training data corresponding to $T = 100$ out of the 180 virus sequences. We train binding and stability prediction models on this data and compute the BROAD antibody sequence by solving an ILP with the T virus sequences in the training subsample. Next, for each training subsample we learn the bi-linear model in Equation 3 and save the coefficients. Then, we solve the MILP 6 to compute the corresponding STRONG antibody for a given α . Given this antibody, we solve ILP 4 to compute T escaping virus sequences corresponding to each of the T training sequences (native). We use CPLEX version 12.51 to solve the (mixed) integer linear programs. Finally, we train a z-score model on the full dataset ($T = 180$). We evaluate the BROAD and the STRONG antibody sequences in terms of the predicted z-score against a) the full 180 virus panel and b) the 100 escaping virus sequences in case of each training subsample. This procedure is outlined in Algorithm 1. As we show in Figure 1

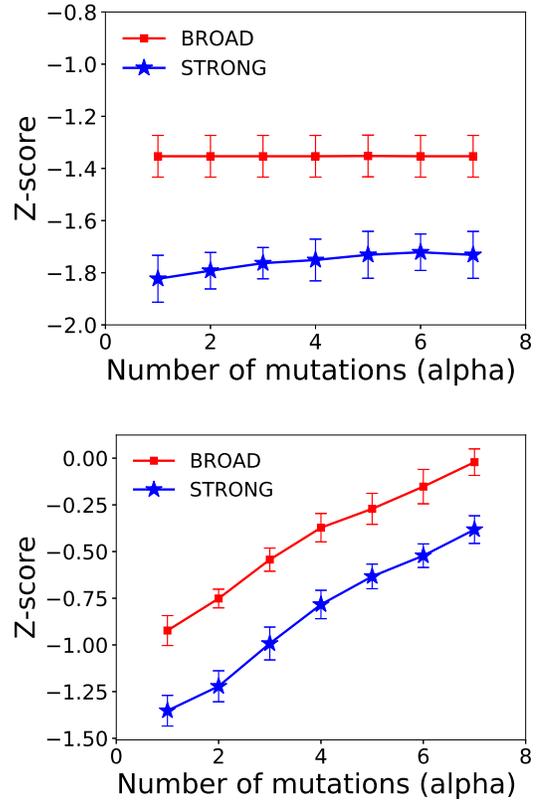


Figure 1: Comparison between STRONG and BROAD in terms of the z-score objective (lower is better): on the full 180 virus panel (left) and the 180 escaping virus set (right).

STRONG is significantly better in minimizing the z-score objective as compared to BROAD.

Algorithm 1 Generating and evaluating STRONG antibody candidates

for each random training set subsample corresponding to $T = 100$ virus sequences **do**

Training Data: $\mathcal{B}(\mathbf{a}, \mathbf{v}), \mathcal{S}(\mathbf{a}, \mathbf{v}), \mathcal{Z}(\mathbf{a}, \mathbf{v})$ corresponding to the T training sequences

Learning: bi-linear model for z-score

Optimization: STRONG antibody \leftarrow MILP 6 solution, escaping set \leftarrow ILP 4 solution

Evaluation: predicted z-score using model trained on the full dataset, and ROSETTA modeling

Finally, we evaluate in terms of the breadth of binding (fraction of viruses in the evaluation panel to which the designed antibody binds) generated using ROSETTA structure modeling. We choose 50 random subsamples of training sets with $T = 30$ virus sequences. Based on binding and stability models trained on the full dataset, we generate the top 10 BROAD candidates. Next, we generate the STRONG antibody for a randomly chosen top BROAD candidate

Table 1: ROSETTA structure modeling results: breadth of binding (%).

Virus Sequences for Evaluation	VRC23	BROAD	STRONG
180 HIV panel	53.3	100	96.1
30 Escaping sequences	43.3	86.7	93.3
30 Training sequences	56.7	100	100

using $\alpha = 5$. We perform ROSETTA structure modeling on these antibody candidates (one BROAD and one STRONG candidate) and the escaping set of 30 virus sequences. For comparison, we also include the native antibody VRC23. We present the ROSETTA computed breadth in each case in Table 1. STRONG significantly outperforms BROAD against the escaping virus panel while it continues to be effective against the training panel.

5.3 Comparison against AAMAS2015

The game-theoretic antibody design approach in [16] uses machine learning guided stochastic local search to compute optimal antibodies. Following the biased random approach in the above research, we generate a set of 350 antibody sequences starting with VRC23. We compute the corresponding average escape costs (number of mutations to escape) with greedy local search starting from the 180 virus panel. We train a binary antibody-virus binding prediction model using an rbf kernel SVM on the full dataset and use this model in the greedy search to evaluate binding. Next, we learn a linear regression model to predict this average escape cost as a function of the antibody sequence. Using these models, we perform 50 independent sequences of local search (400 iterations, random with native bias [16]) to compute 50 optimal antibody candidates.

For comparison, we generate STRONG antibodies corresponding to the above 50 antibodies, with α set to the nearest integer escape cost in MILP 6. Using the z-score model trained on the full dataset, we evaluate these antibodies in Figure 2, on the full 180 panel and the escaping set in each case (from ILP 4). Our proposed approach is significantly better in minimizing the objective (z-score). We also plot the comparison as a function of the local search iterations and observe a similar trend in Figure 3. Note that the z-scores increase with iterations since the average escape cost increases as well.

6 CONCLUSIONS

We proposed an efficient approach for computational antibody design using a Stackelberg game model for the interaction between the antibody and the virus. We formulated the game as a bi-level optimization problem, and proposed a method for solving it by leveraging a bi-linear model predicting binding stability as a function of antibody and virus sequence, combined with integer programming. We show, in particular, that we can compute optimal virus escape using an integer linear program the LP relaxed version of which has integer solutions. Consequently, taking the dual of the associated relaxed LP we obtain an optimization program which can be directly embedded in the optimal antibody design problem, so that the antibody design problem can be solved using mixed-integer

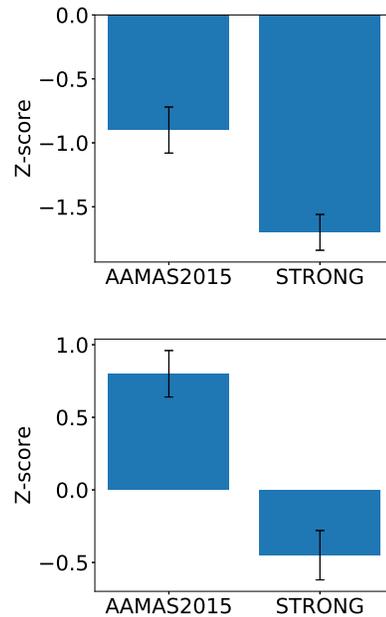


Figure 2: Comparison between STRONG and AAMAS2015 in terms of the z-score objective (lower is better): on the full 180 virus panel (left) and the 180 escaping virus set (right).

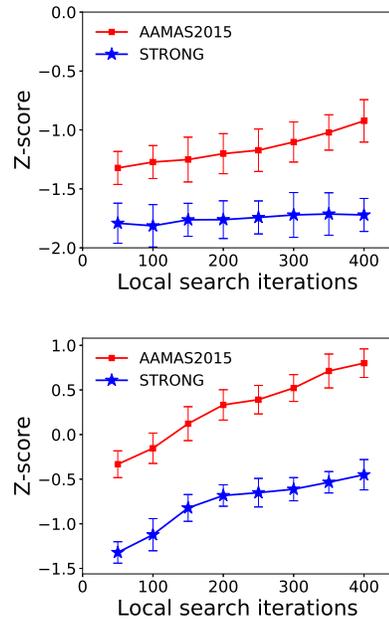


Figure 3: Comparison between STRONG and AAMAS2015 in terms of the z-score objective against search iterations (lower is better): on the full 180 virus panel (left) and the 180 escaping virus set (right).

linear programming. Our experiments show that our approach significantly outperforms both the prior game theoretic alternative, and a state-of-the-art broadly binding antibody design algorithm.

REFERENCES

- [1] Rebecca F Alford, Andrew Leaver-Fay, Jeliuzko R Jeliuzkov, Matthew J O'Meara, Frank P DiMaio, Hahnbeom Park, Maxim V Shapovalov, P Douglas Renfrew, Vikram K Mulligan, Kalli Kappel, et al. 2017. The Rosetta all-atom energy function for macromolecular modeling and design. *Journal of chemical theory and computation* 13, 6 (2017), 3031–3048.
- [2] Chris T. Bauch and David J.D. Earn. 2004. Vaccination and the theory of games. *Proceedings of the National Academy of Sciences* 101, 36 (2004), 13391–13394.
- [3] CDC. 2014. 2014 Ebola Outbreak in West Africa. (2014). <http://www.cdc.gov/vhf/ebola/outbreaks/guinea/>.
- [4] CDC. 2018. Seasonal Influenza, More Information. (2018). <https://www.cdc.gov/flu/about/qa/disease.htm>.
- [5] GB Chapman, M Li, J Vietri, Y Ibuka, D Thomas, H Yoon, and AP. Galvani. 2012. Using game theory to examine incentives in influenza vaccination behavior. *Psychological Science* 23, 9 (2012), 1008–1015.
- [6] Gwo-Yu Chuang, Priyamvada Acharya, Stephen D Schmidt, Yongping Yang, Mark K Louder, Tongqing Zhou, Young Do Kwon, Marie Pancera, Robert T Bailer, Nicole A Doria-Rose, et al. 2013. Residue-level prediction of HIV-1 antibody epitopes based on neutralization of diverse viral strains. *Journal of virology* 87, 18 (2013), 10047–10058.
- [7] Steven A Combs, Samuel L DeLuca, Stephanie H DeLuca, Gordon H Lemmon, David P Nannemann, Elizabeth D Nguyen, Jordan R Willis, Jonathan H Sheehan, and Jens Meiler. 2013. Small-molecule ligand docking into comparative models with Rosetta. *Nature protocols* 8, 7 (2013), 1277.
- [8] José M Cuevas, Ron Geller, Raquel Garijo, José López-Aldeguer, and Rafael Sanjuán. 2015. Extremely high mutation rate of HIV-1 in vivo. *PLoS biology* 13, 9 (2015), e1002251.
- [9] Ivelin S Georgiev, Nicole A Doria-Rose, Tongqing Zhou, Young Do Kwon, Ryan P Staup, Stephanie Moquin, Gwo-Yu Chuang, Mark K Louder, Stephen D Schmidt, Han R Altae-Tran, et al. 2013. Delineating antibody recognition in polyclonal sera from patterns of HIV-1 isolate neutralization. *Science* 340, 6133 (2013), 751–756.
- [10] Jinghe Huang, Byong H. Kang, Marie Pancera, Jeong Hyun Lee, Tommy Tong, Yu Feng, Ivelin S. Georgiev, Gwo-Yu Chuang, Aliaksandr Druz, Nicole A. Doria-Rose, Leo Laub, Kwinten Sliepen, Marit J. van Gils, Alba Torrents de la Pena, Ronald Derking, Per-Johan Klasse, Stephen A. Migueles, Robert T. Bailer, Munir Alam, Pavel Pugach, Barton F. Haynes, Richard T. Wyatt, Rogier W. Sanders, James M. Binley, and Andrew B. Ward. 2014. Broad and potent HIV-1 neutralization by a human antibody that binds the gp41-gp120 interface. *Nature* (2014).
- [11] Manish Jain, James Pita, Milind Tambe, Fernando Ordóñez, Praveen Paruchuri, and Sarit Kraus. 2008. Bayesian stackelberg games and their application for security at Los Angeles international airport. *SIGecom Exch.* 7, Article 10 (June 2008), 3 pages. Issue 2.
- [12] Hetunandan Kamisetty, Bornika Ghosh, Christopher James Langmead, and Chris Bailey-Kellogg. 2015. Learning sequence determinants of protein: Protein interaction specificity with sparse graphical models. *Journal of Computational Biology* 22, 6 (2015), 474–486.
- [13] Jingzhou Liu, Beth F. Kochin, Yonas I. Tekle, and Alison P. Galvani. 2012. Epidemiological game-theory dynamics of chickenpox vaccination in the USA and Israel. *Journal of the Royal Society Interface* 9, 66 (2012), 68–76.
- [14] Joint United Nations Programme on HIV/AIDS et al. 2017. Fact sheet—Latest statistics on the status of the AIDS epidemic. (2017).
- [15] Swetasudha Panda and Yevgeniy Vorobeychik. 2015. Designing Vaccines that Are Robust to Virus Escape. In *AAAI* 4188–4189.
- [16] Swetasudha Panda and Yevgeniy Vorobeychik. 2015. Stackelberg games for vaccine design. In *International Conference on Autonomous Agents and Multiagent Systems*. 1391–1399.
- [17] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in Python. *Journal of machine learning research* 12, Oct (2011), 2825–2830.
- [18] Alexander M Sevy, Tim M Jacobs, James E Crowe Jr, and Jens Meiler. 2015. Design of protein multi-specificity using an independent sequence search reduces the barrier to low energy sequences. *PLoS computational biology* 11, 7 (2015), e1004300.
- [19] Alexander M Sevy, Swetasudha Panda, James E Crowe Jr, Jens Meiler, and Yevgeniy Vorobeychik. 2018. Integrating linear optimization with structural modeling to increase HIV neutralization breadth. *PLoS computational biology* 14, 2 (2018), e1005999.
- [20] UNAIDS. 2016. Fact sheet—Latest statistics on the status of the AIDS epidemic. (2016).