Peer

Machine learning aided multiscale modelling of the HIV-1 infection in the presence of NRTI therapy

Huseyin Tunc¹, Murat Sari² and Seyfullah Kotil^{3,4}

¹ Department of Biostatistics and Medical Informatics, School of Medicine, Bahcesehir University, Istanbul, Turkey

² Mathematics Engineering, Faculty of Science and Letters, Istanbul Technical University, Istanbul, Turkey
 ³ Department of Biophysics, School of Medicine, Bahcesehir University, Istanbul, Turkey

⁴ Department of Molecular Biology and Genetics, Faculty of Arts and Sciences, Bogazici University, Istanbul, Turkey

ABSTRACT

Human Immunodeficiency Virus (HIV) is one of the most common chronic infectious diseases in humans. Extending the expected lifetime of patients depends on the use of optimal antiretroviral therapies. Emergence of the drug-resistant strains can reduce the effectiveness of treatments and lead to Acquired Immunodeficiency Syndrome (AIDS), even with antiretroviral therapy. Investigating the genotype-phenotype relationship is a crucial process for optimizing the therapy protocols of the patients. Here, a mathematical modelling framework is proposed to address the impact of existing mutations, timing of initiation, and adherence levels of nucleotide reverse transcriptase inhibitors (NRTIs) on the evolutionary dynamics of the virus strains. For the first time, the existing Stanford HIV drug resistance data have been combined with a multi-strain within-host ordinary differential equation (ODE) model to track the dynamics of the most common NRTI-resistant strains. Overall, the D4T-3TC, D4T-AZT and TDF-D4T drug combinations have been shown to provide higher success rates in preventing treatment failure and further drug resistance. The results are in line with the genotypephenotype data and pharmacokinetic parameters of the NRTI inhibitors. Moreover, we show that the undetectable mutant strains at the diagnosis have a significant effect on the success/failure rates of the NRTI treatments. Predictions on undetectable strains through our multi-strain within-host model yielded the possible role of viral evolution on the treatment outcomes. It has been recognized that the improvement of multi-scale models can contribute to the understanding of the evolutionary dynamics, and treatment options, and potentially increase the reliability of genotype-phenotype models.

Subjects Computational Biology, Mathematical Biology, HIV, Statistics **Keywords** AIDS, HIV infection, Machine learning, NRTI therapy, Mathematical models

INTRODUCTION

Antiretroviral drug resistance is one of the main barriers to therapy success for HIVpositive patients. According to the WHO, the HIV drug resistance report 2021, 10% and 40% of adults are affected by drug-resistant strains (DRS) for naive and treated patients, respectively. In addition, 50% of newly diagnosed infants were exposed to the DRS. The

Submitted 17 June 2022 Accepted 19 February 2023 Published 31 March 2023

Corresponding author Seyfullah Kotil, enesseyfullah.kotil@med.bau.edu.tr, enesseyfullah.kotil@boun.edu.tr

Academic editor Yuriy Orlov

Additional Information and Declarations can be found on page 0

DOI 10.7717/peerj.15033

Copyright 2023 Tunc et al.

Distributed under Creative Commons CC-BY 4.0

OPEN ACCESS

DRS can be acquired with nonadherence to the therapy protocols, or patients can directly be infected with DRSs (*Blower et al., 2001*). Both scenarios yield life-long persistence of the DRS and need to be carefully tracked.

Quantitative evaluation of HIV drug resistance has been carried out with the use of phenosense assays by finding the fold-change of IC_{50} values (the amount of concentration to inhibit 50% of virion) between drug-resistant and wild-type strains (*Zhang et al., 2005; Pham et al., 2018; Feng et al., 2016*). Data modelling frameworks have been used to construct general mathematical relations between genotype and phenotype information (*Tarasova et al., 2018; Steiner, Gibson & Crandall, 2020; Shah et al., 2020; Tarasova et al., 2023; Lagunin et al., 2023*). These mathematical models aim to generalize the given data by means of encoding the amino acid sequence of target enzymes (*Rhee, Taylor & Fessel, 2010*). One of the main contributions of the current study is to explore how these models can be embedded into a within-host model to simulate the evolutionary dynamics of HIV strains. In particular, we simulate thousands of clinically relevant HIV mutant strains provided by *Rhee, Taylor & Fessel (2010*).

For forecasting the viral dynamics of HIV, various within-host models have been presented in ordinary differential equation (ODE) forms in the presence/absence of resistant strains and antiretroviral therapy (Perelson & Nelson, 1999; Dixit & Perelson, 2004; Rong, Feng & Perelson, 2007; Hadjiandreou, Conejeros & Vassiliadis, 2007; Sutimin et al., 2017; Wu & Zhao, 2020; Chen, Teng & Zhang, 2021). Perelson & Nelson (1999) proposed HIV within-host models consisting of CD4+ T cells (T), infected CD4+ T cells (T^*), macrophage cells (M), infected macrophage cells (M^*) and virions (V) in the presence and absence of antiretroviral therapy. Dixit & Perelson (2004) proposed $T - T^* - V$ model by considering time-dependent intracellular efficiency of reverse transcriptase and protease inhibitors (RTIs and PIs). Rong, Feng & Perelson (2007) derived two-strain extension of the within-host model given in the literature (Perelson & Nelson, 1999; Dixit & Perelson, 2004) with antiretroviral therapy. Hadjiandreou, Conejeros & Vassiliadis (2007) revised the $T - T^* - M - M^* - V$ model proposed by *Perelson & Nelson (1999)* by adding homeostatic cell proliferation terms to capture long time behaviour of the HIV dynamics. Sutimin et al. (2017) modeled the within-host HIV dynamics with target Langerhans and CD4+ T cells and investigated the time-dependent efficiency of RTIs and PIs with various scenarios. Wu & Zhao (2020) derived two-strain within-host model including the age of infection detail represented by the system of integro-differential equations. They mathematically formulated the competition between the drug-sensitive and drug-resistance strains with respect to model parameters. Chen, Teng & Zhang (2021) included saturated incidence and distributed infection delays into the standard two-strain $T - T^* - V$ model and investigated the effects of those novel incidence terms on the long-time behaviour of the dynamics. Additionally, the effect of drug adherence on the virological failure of ARTs (Rosenbloom et al., 2012), the effect of time-dependent drug efficiencies on ART response (Rong, Feng & Perelson, 2007; Vaidya & Rong, 2017), competition between susceptible and resistant strains in the viral dynamics (Ball, Gilchrist & Coombs, 2007; Lythgoe, Pellis & Fraser, 2013), the role of latently infected CD4+ T cell reservoirs on the evolution of strains (Doekes, Fraser & Lythgoe, 2017), comparison of entry inhibitors with the RTIs and PIs according to viral

resistance (*Alshorman, Al-hosainat & Jackson, 2022*), investigation of optimal timing for ART *Rouzine (2022)* have been proposed through within-host models. The proposed mathematical models assume the co-existence of susceptible and resistant strains and generally investigate the response to antiretroviral therapy (ART). The current study addresses similar questions with a novel multiscale model based on Stanford HIV Drug Resistance data and machine learning models.

For the first time, we combined the experimental drug resistance data of nucleotidereverse transcriptase inhibitors (NRTI) available in the Stanford HIV drug resistance database (*Rhee et al., 2003*) with a within-host model of HIV infection to observe the dynamics of the viral strains under different scenarios. Our multiscale model brings together three pieces of information: IC_{50} values for each mutant with machine learning models, within blood dynamics for NRTIs, and CD4+ T cells and macrophage cells for primary targets of virions. For different mutant compositions, we aim to investigate the emergence of treatment failure for different initiation timing (up to one year) and adherence level of NRTI therapies (21 different combinations). Here we rank the inhibitory capabilities of the NRTI combinations in the presence of various viral strains and ongoing viral evolution. Our results add to the predictions of the Stanford HIV drug resistance database, which identifies the best drug by selecting the one that has the lowest IC_{50} for a given mutant. But that model is a static model that cannot incorporate the effects of new mutants that can be generated through time which is accounted for in our model.

MATERIALS AND METHODS

Within-host model with wild-type virus

In this part, we have inspired from the earlier studies on the within-host HIV infection model (Hadjiandreou, Conejeros & Vassiliadis, 2007; Hernandez-Vargas, 2019; Hernandez-Vargas & Middleton, 2013). We assume that the primary reservoirs for HIV infection are: CD4+T cells and macrophages denoted by T(t) and M(t) (Hernandez-Vargas, 2019; Hernandez-Vargas & Middleton, 2013). The long-living macrophage cells cause the persistence of virions over the years (Orenstein, 2001; Herbein & Varin, 2010). Macrophage cells contribute to the depletion of healthy CD4 + T cells in advanced HIV infection (*Crowe*, 1995). Within-host modelling of HIV infection without considering the macrophage reservoirs yielded less reliable dynamics, such as the models that never result in the AIDS phase (Rong, Feng & Perelson, 2007). We denote the HIV infected CD4+ T cells and macrophages by $T^*(t)$ and $M^*(t)$. Lastly, the number of free wild-type virions in the host is denoted by the function V(t). By considering model assumptions like homeostatic cell proliferation terms (s_T, s_M) , bilinear incidence terms $(k_T TV, k_M TM)$, natural deaths of cells and virions ($\delta_T T$, $\delta_M M$, $\delta_{T^*} T^*$, $\delta_{M^*} M^*$, $\delta_V V$), viral replication terms $(p_T T^*, p_M M^*)$ and the Michaelis–Menten type proliferation terms $\left(\frac{\rho_T V}{c_T + V}T, \frac{\rho_M V}{c_M + V}M\right)$, we express the one strain within-host model with the following system of ordinary differential equations (Hernandez-Vargas, 2019; Hernandez-Vargas & Middleton, 2013)

$$\frac{dT}{dt} = s_T - k_T T V - \delta_T T + \frac{\rho_T V}{c_T + V} T$$

$$\frac{dT^*}{dt} = k_T T V - \delta_{T^*} T^*$$

$$\frac{dM}{dt} = s_M - k_M M V - \delta_M M + \frac{\rho_M V}{c_M + V} M$$
(1)
$$\frac{dM^*}{dt} = k_M M V - \delta_{M^*} M^*$$

$$\frac{dV}{dt} = p_T T^* + p_M M^* - \delta_V V$$

where initial conditions are considered as $T(0) = T_0$, $T^*(0) = T_0^*$, $M(0) = M_0$, $M^*(0) = M_0^*$ and $V(0) = V_0$. Further details of the model (1) can be seen in the study of *Hernandez-Vargas & Middleton (2013)*. In the following section, we expand the model Eq. (1) to include both susceptible and resistant multiple strains as well as NRTI therapy.

Multi strain within-host model with NRTI therapy

The ARTs include at least one of the NRTIs that aim to block the activation of the reverse transcriptase enzyme. Effective treatment of HIV-positive patients with NRTIs saves millions of lives worldwide (*Tressler & Godfrey, 2012*). However, the error-prone structure of the HIV replication yields resistant strains over the years, and these strains are known to be a primary barrier to preventing AIDS (*Kuritzkes, 2011*). Our multiscale within-host model includes three main steps: constructing machine learning models to generalize isolate-fold change data for NRTIs, a model for dealing with NRTI action in blood, and finally, a within-host model with multi-strains and NRTI therapy.

An artificial neural network model for isolate-fold change relation

There exists various genotype-phenotype experiment data, including the fold change values of *IC*₅₀ (the required drug concentration to inhibit 50% of virions) for various reverse transcriptase inhibitors in the presence of susceptible and resistant isolates (*Rhee et al., 2005*). The most used genotype-phenotype data is the Stanford HIV drug resistance database (https://hivdb.stanford.edu/). We use filtered genotype-phenotype data of reverse transcriptase inhibitors available in this database and are widely used for various machine learning algorithms (*Amamuddy, Bishop & Bishop, 2017; Masso & Vaisman, 2013*). By regulating the data for each NRTI, 1,224 unique mutations were observed for the reverse transcriptase enzyme. In this filtered dataset, 1,662 isolates for epivir (3TC), 1,597 isolates for abacavir (ABC), 1,683 isolates for zidovudine (AZT), 1,693 isolates for stavudin (D4T), 1,693 isolates for didanosine (DDI) and 1,354 isolates for tenofovir (TDF) have been analyzed for NRTI susceptibility. The dataset includes 1,206, 1,136, 1,220, 1,223, 1,223, and 1,119 unique mutations for 3TC, ABC, AZT, D4T, DDI, and TDF, respectively.

Here, we apply the binary barcoding technique (*Rhee, Taylor & Fessel, 2010*) to represent the isolates occurring in the dataset. Hence, 1,224-dimensional input vectors of 0s and 1s are created by considering the existence of unique mutations in the isolates. Let us denote our complete mutation set as $M = \{m_1, m_2, ..., m_{1224}\}$ where m_i is an NRTI specified

mutation pattern. We define the binary representation of isolate *j* as $I_j = [a_1, a_2, ..., a_{1224}]$ with

$$a_k = \begin{cases} 1, & \text{if } m_k \in I_j \\ 0, & \text{otherwise.} \end{cases}$$

We construct six artificial neural networks (*ANN*) models to predict logarithmic foldchange values in the presence of any isolates related to each NRTI therapy by using the Machine Learning and Deep Learning toolbox of the MATLAB 2022a program (https://www.mathworks.com/). The *ANN* architectures include 1,224-dimensional input, five hidden layer neurons, and one output neuron with hyperbolic tangent-sigmoid and linear activation functions. The model selection process is explained with detailed quantitative observations in Table S1. The scaled conjugate gradient algorithm with MATLAB built-in function "trainscg" has been used in the training process over GPU. Let us denote our model as a function that maps isolate vectors to the fold changes as

Fold $Change = ANN_X(I)$

where $I \in \{0, 1\}^{1 \times 1224}$ and X is a specified inhibitor ($X \in \{3\text{TC}, \text{ABC}, \text{AZT}, \text{D4T}, \text{DDI}, \text{TDF}\}$). To overcome possible overfitting, we have implemented an ensemble learning process. For each inhibitor, the 50 ×100 model has been trained with random training, validation, and test set (80%, 10% and 10%). A model is chosen from every 100 models that yield the minimum mean square error for the test set of the corresponding inhibitor data. Hence, 50 optimal models are selected out of 5,000 models for each NRTI inhibitor, and the final model is calculated as the average of these models.

The prediction performance of six $ANN_X(I)$ models with linear correlation coefficient (R) and mean square error (MSE) values are presented in Fig. 1. According to the figure, $ANN_{X}(I)$ models yield accurate predictions with high R and low MSE scores. Mean MSE value of $ANN_X(I)$ models have been obtained as 0.0453 with 95% CI [0.0005–0.0901]. Similarly, the mean R value of the models has been calculated as 0.9093 with 95% CI [0.8677-0.9509]. To observe how six $ANN_X(I)$ models classify resistant and susceptible strains, we convert our regression models into classification models by labeling the data as resistant (Fold Change \geq 3) and susceptible (Fold Change < 3). The receiving operating curves (ROC) corresponding to the six ANN models and the area under the curve (AUC) values are presented in Fig. S1. According to the classification results, we get the mean AUC score as 0.9649 with 95% CI [0.9423-0.9875]. Additionally, to see why such a nonlinear model is needed to map the genotype data into the phenotype output, we also perform multiple linear regression (MLR) analysis (with 20% holdout data) for data of six NRTIs. The regression and classification performance of the MLR models are shown in Figs. S2–S3. A fair comparison between the ANN and MLR models in terms of the MSE, R, and AUC values is given in Table S2. According to the table, even classification performance of the models is almost the same, the ANN models give much more accurate estimations in regression. Since better regression performance is more desirable for our further modelling framework, the ANN models are assumed to be our baseline models for predicting the resistance profiles of given viral strains.





Full-size DOI: 10.7717/peerj.15033/fig-1

Modelling the time-dependent drug efficacy

Modelling the efficacy of antiretrovirals using the plasma drug concentrations can be seen in various studies in the literature (*Dixit & Perelson, 2004*; *Rong, Feng & Perelson,* 2007; *Rosenbloom et al., 2012*). *Rosenbloom et al. (2012*) modeled the time-dependent drug efficiency in plasma by considering the exponential decay of plasma drug concentration after the instantaneous peak. Here we use the time-dependent drug efficacy model described by *Dixit & Perelson (2004)* and *Rong, Feng & Perelson (2007)*, considering the pharmacokinetic parameters of drugs in the blood. *Dixit & Perelson (2004)* considered the phosphorylated concentration of the tenofovir (TDF) in the cells. Since the time-drug efficiency functions obtained by taking into account blood concentration and phosphorylated within cell concentration of drugs follow a very similar trend, here we assume the blood concentration of the drugs (see Fig. 1 of *Dixit & Perelson (2004)*). Additionally, the non-availability of phosphorylation reaction parameters for the remaining five inhibitors 3TC, ABC, AZT, D4T, and DDI have encouraged us to consider the blood concentration of the drugs only.

Let $\varepsilon_X^Y(t)$ denotes the time-dependent efficacy of drug X in the presence of strain (isolate) Y. The instantaneous efficacy can be approximated as *Dixit & Perelson* (2004)

$$\varepsilon_X^Y(t) = \frac{C_b^X(t)}{(IC_{50})_X^Y + C_b^X(t)}$$
(2)

where $C_b^X(t)$ denotes the within blood concentration of drug X and $(IC_{50})_X^Y$ denotes the required concentration of drug X to inhibit the 50% of strain Y. According to our

isolate-fold change ANN model, Eq. (2) can be rewritten as

$$\varepsilon_X^Y(t) = \frac{C_b^X(t)}{ANN_X(Y)(IC_{50})_X^{WT} + C_b^X(t)}$$
(3)

where $(IC_{50})_X^{WT}$ denotes the required concentration of drug *X* to inhibit the 50% wild type virus. Thus, to completely describe $\varepsilon_X^Y(t)$, we should model $C_b^X(t)$. According to *Dixit & Perelson (2004)*, the concentration of a drug in the blood can be expressed as

$$C_{b}(t) = \frac{FDk_{a}e^{-k_{e}t}}{V_{d}(k_{e}-k_{a})(e^{k_{a}I_{d}}-1)} \left[1 - e^{(k_{e}-k_{a})t}\left(1 - e^{N_{d}k_{a}I_{d}}\right) + \frac{\left(e^{k_{e}I_{d}} - e^{k_{a}I_{d}}\right)\left(e^{(N_{d}-1)k_{e}I_{d}}-1\right)}{e^{k_{e}I_{d}}-1} - e^{((N_{d}-1)k_{e}+k_{a})I_{d}}\right]$$
(4)

where *F* is the bioavailability of the drug, *D* is the mass of the drug administered in one dose, I_d is the dosing interval, N_d is the number of doses up to time *t*, V_d is the volume of distribution, k_a and k_e are pharmacokinetic parameters. The drug-specific parameters k_a , k_e , *D*, I_d and *F* occurred in Eq. (4) and IC_{50} values for 3TC, ABC, AZT, D4T, DDI, and TDF according to the equations given by *Dixit & Perelson (2004)* are evaluated and presented in Table 1. Detailed explanations of the derivation of these parameters are given in the Supplementary Information.

A multi-strain within-host model

This part of the study combines our investigations into a unique multi-strain within-host model. To reduce the cost of the simulations, we assume the main NRTI-related mutations 115F, 151M, 184I, 184V, 210W, 215F, 215Y, 41L, 65N, 65R, 67N, 69D, 70E, 70G, 70R, 74I and 74V according to the study of *Rhee et al.* (2005). These 17 mutations yield 131,071 unique strains having all possible mutations. Thus, by considering wild-type and mutant strains, we have total N = 131,072 strains. Our multi-strain within-host model with time-dependent NRTI therapy can be derived from one strain model (1) as follows

$$\frac{dT}{dt} = s_T - k_T T \sum_{i=1}^{N} (1 - c_i)(1 - \varepsilon_X^i(t))V_i - \delta_T T + \frac{\rho_T \sum_{i=1}^{N} V_i}{c_T + \sum_{i=1}^{N} V_i} T$$

$$\frac{dT_i^*}{dt} = k_T (1 - c_i)(1 - \varepsilon_X^i(t))TV_i - \delta_{T^*} T_i^*$$

$$\frac{dM}{dt} = s_M - k_{MM} \sum_{i=1}^{N} (1 - c_i)(1 - \varepsilon_X^i(t))V_i - \delta_M M + \frac{\rho_M \sum_{i=1}^{N} V_i}{c_M + \sum_{i=1}^{N} V_i} M$$
(5)
$$\frac{dM_i^*}{dt} = k_M (1 - c_i)(1 - \varepsilon_X^i(t))MV_i - \delta_{M^*} M^*$$

$$\frac{dV_i}{dt} = p_T T_i^* + p_M M_i^* - \delta_V V_i$$

Table 1 Drug specific parameters for time-dependent drug efficiency equation Eq. (4).							
Parameter/Drug	3TC	ABC	AZT	D4T	DDI	TDF	
$IC_{50}(\times 10^{-5} \text{ mg/ml})$	3.97	132.64	1.87	4.25	113.11	16.24	
D(mg)	300	300	300	40	400	300	
I_d (day)	1	0.5	0.5	0.5	1	1	
F	0.86	0.83	0.64	0.86	0.42	0.39	
k_a	27.98	51.07	37.42	54.29	32.34	8.36	
k_e	3.44	8.52	14.25	7.84	47.30	16.58	
V_d (ml)	91,000	60,200	112,000	46,000	54,000	87,500	

where i = 1, 2, ..., N = 131,072, $T_i^*(t)$ and $M_i^*(t)$ denote the number of CD4 + T cells and macrophage cells infected by strain *i* and $V_i(t)$ represents the number of virions having *i* th genotype. In the multi-strain within-host model (5), $\varepsilon_X^i(t)$ denotes the time-dependent efficacy of the inhibitor X on the strain *i* and $0 \le c_i \le 1$ represents the fitness costs of mutant strains with $c_1 = 1$ for the wild type of strain. The lack of enough experimental results on these fitness values compelled us to use the mean fitness cost values of mutations 41L, 67N, 70R, 184V, 210W, 215D, 215S, and 219Q estimated by *Kühnert et al. (2018)* as 0.2232, 0.3181, 0.3863, 0.5899, 0.3091, 0.0981, 0.1664 and 0.3207, respectively. According to these data, we assume that $c_i = 0.3015$ for mutant strains $i \ge 2$. A schematic illustration of the multi-strain within-host model (5) is given in Fig. 2. Parameter values of multi-strain within-host model (5) with corresponding references can be seen in Table 2.

The within-host model (5) ignores the role of latently infected CD4+ T cells. The main role of latently infected CD4+ T cells is the viral rebound after poor adherence to the given therapy (Chun et al., 2000), and these cells are almost three percent of all CD4+ T cells (Hadjiandreou, Conejeros & Wilson, 2009). Since model (5) is continuous over time and hence the emerged viral strains are not completely eradicated in the viral suppression phase, the persistence of HIV-1 virions is automatically ensured, and poor adherence in model (5) provides viral rebound. Thus, ignoring the latently infected CD4+ T cells in model (5) does not considerably affect our modelling framework. As indicated in the study of Chun et al. (2000), latently infected CD4+ T cells are not the only reason for the rebound of plasma viremia after discontinuation of the ART. The literature (Alexaki, Liu & Wigdahl, 2008; Hendricks et al., 2021; Kruize & Kootstra, 2019) show that the macrophage cells are of particular importance in HIV-1 persistence, and this is why model (5) considers this observation like some existing studies (Hadjiandreou, Conejeros & Vassiliadis, 2007; Hadjiandreou, Conejeros & Wilson, 2009; Hernandez-Vargas, 2019; Hernandez-Vargas & *Middleton*, 2013). Consideration of the role of the macrophage cells yields slow progression to the AIDS phase for untreated patients and improves the reliability of model outcomes (Hernandez-Vargas, 2019).

To model the effect of mutations, we do not explicitly include the mutation matrix in the ODE system (5); instead, we address the transition between mutations and strains at the end of each time step by generating Poisson random numbers (*Rosenbloom et al., 2012*). Let us assume time step n (t = n day), $(T_i^*)_n = T_i^*(n)$ and $(M_i^*)_n = M_i^*(n)$. The mutation

Peer.



Figure 2 Illustration of the core parts of multi-strain within-host model (5) with NRTI therapy. Model (5) assumes the healthy CD4 + T cells (T(t)) and macrophage cells (M(t)) as the main targets of the viral strains $(V_i(t))$. T(t) and M(t) increase with both homeostatic cell proliferation and cell proliferation due to the increasing viral load. Viral strains infect both CD4 + T cells and macrophage cells and then those healthy cells become infected CD4 + T cells $(T_i^*(t))$ and macrophage cells $(M_i^*(t))$. $T_i^*(t)$ and $M_i^*(t)$ compartments produce mature viral strains $V_i(t)$ with some constant rates. All compartments have natural death or clearance with some constant rates. NRTIs block the infection mechanism of the viral strains in healthy cells. The efficiency of the NRTIs is estimated through pharmacokinetic Eq. (3) and the pretrained artificial neural network models that map the genotype data to fold-change values of the IC_{50} 's with respect to the wild type virion.

Full-size DOI: 10.7717/peerj.15033/fig-2

matrix of our system is denoted by MT and defined as

$$MT_{ij} = \begin{cases} 1, & \text{if strain i can take a mutation to become strain j} \\ 0, & \text{otherwise} \end{cases}$$
(6)

For the infected CD4 T cells $(T_i^*)_n$ and infected macrophage cells $(M_i^*)_n$, we calculate the number of new infected ones in one day period as $\Delta(T_i^*)_n$ and $\Delta(M_i^*)_n$ without taking into account the death of these newly infected cells. For each i = 1, 2, ..., N, poissrnd $(\mu \Delta(T_i^*)_n)$

Parameter	Value	Unit	Parameter variation
s_T	10^{4}	$ml^{-1}d^{-1}$	$7 \times 10^3 - 2 \times 10^4$
s_M	150	$ml^{-1}d^{-1}$	100 - 300
k_T	4.5714×10^{-8}	mld^{-1}	$3.2 \times 10^{-8} - 10^{-7}$
k_M	4.3333×10^{-11}	mld^{-1}	$1.73 \times 10^{-11} - 1.3 \times 10^{-9}$
Pт	38	d^{-1}	30.4 - 114
p_M	35	d^{-1}	22 - 132
δ_T	0.01	d^{-1}	0.001 - 0.017
δ_{T^*}	0.4	d^{-1}	0.1 - 0.45
δ_M	0.001	d^{-1}	$10^{-4} - 1.4 \times 10^{-3}$
δ_{M^*}	0.001	d^{-1}	$10^{-4} - 1.2 \times 10^{-3}$
δ_V	2.4	d^{-1}	0.96 - 2.64
$ ho_T$	0.01	d^{-1}	_
$ ho_M$	0.003	d^{-1}	_
c_T	3×10^{5}	ml^{-1}	_
c_M	2.2×10^{5}	ml^{-1}	-

Table 2Parameter values, units and parameter intervals of the within-host models and taken from theliterature (Hernandez-Vargas & Middleton, 2013; Hernandez-Vargas, 2019).

and $poissrnd(\mu \Delta (T_i^*)_n)$ number of infected cells are randomly transmitted from strain *i* to strain *j* according to the mutation matrix MT_{ij} where function poissrnd(x) generates Poisson random number with mean *x* and $\mu = 3 \times 10^{-5}$ denoting the mutation rate (*Rosenbloom et al., 2012*). Note that the mutation rate for each point mutation is unique for the corresponding amino acid change, but we assume a fixed average mutation rate $\mu = 3 \times 10^{-5}$ as stated by *Rosenbloom et al. (2012)*. Since NRTI-related mutation rates have low variance value (*Rosenbloom et al., 2012*) and we have so many viral strains to track, we use overall mutation rate $\mu = 3 \times 10^{-5}$. Parameter values of models (1) and (5) are presented with their references in Table 2.

Model (5) can also include dual therapy of NRTIs X and Y by modifying the therapyrelated time-dependent infection coefficients for CD4 + T cells and macrophage cells $\beta_i^{T/M}(t) = k_{T/M}(1-c_i)(1-\varepsilon_X^i(t))$ with the use of Bliss independence of drug actions as *Jilek et al. (2012)*

$$\beta_i^{T/M}(\varepsilon_X^i(t), \varepsilon_Y^i(t)) = k_{T/M}(1 - c_i) \left(1 - \varepsilon_X^i(t)\right) \left(1 - \varepsilon_Y^i(t)\right)$$
(7)

or Loewe additivity of drug actions (Jilek et al., 2012)

$$\beta_i^{T/M}(\varepsilon_X^i(t), \varepsilon_Y^i(t)) = k_{T/M}(1 - c_i) \frac{1}{\frac{\varepsilon_X^i(t)}{1 - \varepsilon_X^i(t)} + \frac{\varepsilon_Y^i(t)}{1 - \varepsilon_Y^i(t)} + 1}.$$
(8)

Bliss independence assumes independent actions of combined drugs, and Loewe additivity assumes the competition for the same binding site. According to *Jilek et al.* (2012), all combinations except AZT-D4T and DDI-TDF obey the Bliss independence rule, and these two combinations obey the Loewe additivity rule. Note that, since we assume $k_M \approx k_T/1000$ and $\beta_i^T(t) \approx \beta_i^M(t)/1000$ according to the *Hernandez-Vargas* (2019) and *Hernandez-Vargas & Middleton* (2013) (see Table 2), we prefer to use the notation β_i for β_i^T throughout the following parts. Whenever β_i values are quantitatively mentioned in the results section, these values correspond to the β_i^T .

Note that even though we describe our model parameters for 1 ml of blood in Table 2 as widely assumed in the literature (*Hadjiandreou, Conejeros & Vassiliadis, 2007; Hernandez-Vargas, 2019; Hernandez-Vargas & Middleton, 2013*), we simulate the viral dynamics in the host plasma (3,000 ml; *Rosenbloom et al., 2012*) to catch more viral diversity. We assume that the only reservoir of HIV virions is the plasma, which is the major one (*Valcour et al., 2012*), even if there exist other reservoirs like lymph nodes or cerebrospinal fluid (CSF) (*Valcour et al., 2012; Haase, 1999*). Since the instantaneous drug efficiency rates are ($\varepsilon_X^Y(t)$) in non-dimensionless form, we can easily simulate the dynamics in the host plasma by converting the volume-dependent model parameters given in Table 2. For example, by considering 3L host plasma (*Rosenbloom et al., 2012*), the infectivity parameter $k_T = 4.5714 \times 10^{-8}$ ml/day equivalently becomes $k_T = \frac{4.5714 \times 10^{-8}}{3000}$ plasma/day = 1.5238 × 10⁻¹¹ plasma/day.

RESULTS

This section provides the simulation results of the multi-strain within-host model (5), starting with various viral strains. The effects of adherence levels and initiation timing of NRTI therapies on the progression of viral dynamics are investigated. This section includes four subsections in which we propose the statistics of the infection rates, details of model simulations, the quantitative measure for the therapy success, and the simulation results for various cases.

Statistics of infection rates

Before running the simulations to observe the failure/success distribution of each NRTI combination, we may predict the best possible therapy protocol through our pretrained machine learning model and the pharmacokinetic properties of the inhibitors. Obviously, as we infer from our model (5) and drug-specific time-dependent infection rate $\beta_i(\varepsilon_X^i(t), \varepsilon_Y^i(t))$ (7)–(8), each viral strain has its infection rate and aims to be dominant by infecting more healthy cells. Since evaluation of $\beta_i(\varepsilon_X^i(t), \varepsilon_Y^i(t))$ is straightforward through Eqs. (7)–(8) and (3), we may have some prior estimates for the selection of the best therapy protocol. Distribution of 131,071 $\beta_i(\varepsilon_X^i, \varepsilon_Y^i) = \int_0^1 \beta_i(\varepsilon_X^i(t), \varepsilon_Y^i(t)) dt$ values in the presence of 21 different mono and dual NRTI therapies are illustrated in Fig. 3. Descriptive statistic values of $\beta_i(\varepsilon_X^i, \varepsilon_Y^i)$ values for all combinations are presented in Table 3.

Figure 3 and Table 3 show that the probability distributions are almost uniform and $\beta_i(\varepsilon_X^i, \varepsilon_Y^i)$ values have considerable diversity and standard deviations among the viral strains. Hence, this observation means that even having point mutations can change the infection rates considerably and thus may lead to a need for more perfect adherence levels to the given therapy. Additionally, Fig. 3 implies that the initial viral strain of the patient plays a critical role in the progression of HIV dynamics. According to Table 3, NRTI therapy combinations yield 38.4% and 78% decrease in infection rate on average (among all therapies) (95% CI [36.2%–40.7%] and [69.7%–86.3%]) for the worst and best case scenario (having most and least resistant initial strain), respectively.

Table 3 ranks the possible NRTI combinations in terms of the resistance scores but ignores the side effects and cost-effectiveness. Various side-effects of NRTIs linked with



Figure 3 Probability distributions of infection rate (β_i) values of various viral strains in the presence of NRTI therapy combinations. (β_i) values are calculated with Eqs. (7)–(8) depending on the drug pairs. (β_i) values are effected by pharmacokinetic parameters, IC_{50} values for the viral strains, baseline infection rate $k_T = 4.5714 \times 10^8$ and the fixed viral fitness value ($c_i = 0.3015$) of the viral strains.

Full-size DOI: 10.7717/peerj.15033/fig-3

Drugs	Mean	Min	Max	Std	Median	Mode	Q_1	Q_3
D4T-3TC	1.290	0.160	2.069	0.363	1.295	0.16	1.029	1.576
D4T-AZT	1.370	0.427	2.358	0.442	1.368	0.427	1.021	1.722
TDF-D4T	1.403	0.604	2.373	0.371	1.382	0.604	1.109	1.679
D4T	1.442	0.697	2.319	0.351	1.426	0.697	1.157	1.72
AZT-3TC	1.473	0.154	2.212	0.462	1.531	0.154	1.109	1.877
D4T-ABC	1.525	0.695	2.405	0.374	1.514	0.695	1.221	1.826
DDI-D4T	1.592	0.765	2.466	0.382	1.581	0.765	1.279	1.903
TDF-AZT	1.627	0.474	2.523	0.506	1.683	0.474	1.226	2.056
AZT-ABC	1.755	0.551	2.529	0.503	1.844	0.551	1.363	2.194
AZT	1.775	0.554	2.564	0.513	1.858	0.554	1.373	2.223
DDI-AZT	1.834	0.562	2.602	0.533	1.933	0.562	1.412	2.307
TDF-3TC	1.884	0.3	2.265	0.307	1.952	0.3	1.835	2.065
3TC	1.965	0.29	2.173	0.339	2.114	0.29	2.000	2.133
ABC-3TC	2.030	0.274	2.318	0.359	2.173	0.274	2.025	2.225
DDI-3TC	2.155	0.323	2.373	0.356	2.305	0.323	2.201	2.325
TDF-ABC	2.172	1.508	2.588	0.181	2.176	1.508	2.038	2.314
TDF	2.299	1.889	2.665	0.172	2.3	1.889	2.163	2.438
TDF-DDI	2.323	1.917	2.667	0.164	2.327	1.917	2.194	2.457
ABC	2.459	1.733	2.698	0.126	2.485	1.733	2.404	2.544
DDI-ABC	2.546	1.869	2.726	0.106	2.57	1.869	2.505	2.617
DDI	2.746	2.675	2.78	0.02	2.747	2.675	2.732	2.762

Table 3	Descriptive statistics (x 10 ⁻⁸) of infecti	on rate β _i value	es for all possibl	e mono and dua	l NRTI
therapie	s.					

mitochondrial toxicity (*Holec et al., 2018*). We present the possible side-effects of the existing NRTIs in Table S3, and a detailed review can be found in the study of *Montessori et al. (2004)*. The cost-effectiveness of NRTI therapies is essential to maximize the expected

survival times of the patients with minimized costs. Various mathematical models are available that compare treatments for cost-effectiveness, and a detailed review of *Mauskopf* (2013) provides various essential results. Most of the models described in their study ignore the effect of drug resistance. Drug resistance is a crucial contributor to the expected costs. This study is only interested in the impact of drug resistance on the NRTI therapy outcomes, and we both ignore side effects and cost-effectiveness.

Details of model simulations

In our simulations, we investigate the effect of the type of NRTI therapy, timing of the NRTI therapy, and adherence to the provided therapy on CD4+ T cell counts of the patients. All possible 21 mono and dual NRTI combinations of six inhibitors have been included in the simulations by considering their independent or additive actions. The initiation time of the NRTI therapy is considered within the first year after the patient became infected and denoted by τ . The adherence level of a patient to the provided therapy protocol is assigned to a real number α between 0 and 1, representing nonadherence to full adherence levels. After initiating the treatment with adherence level α in a day of the simulation, the patient takes drug(s) with probability α according to the parameters given in Table 1. Initial viral load, CD4 + T cell count, and macrophage cell count in the simulations are considered as 1 virion/ml, 10⁶ cell/ml and 150 cell/ml, respectively (*Hernandez-Vargas, 2019*).

It is assumed that the patient is infected with one type of mutant strain with one to five-point mutations on the reverse transcriptase enzymes. In this way, five groups are constructed to include five different strains. These viral strains have been determined according to the frequency of presence in the Stanford HIV drug resistance database. These initial viral strains are denoted by G_{ij} where i = 1, 2, 3, 4, 5 denotes the number of the point mutations in the strain and i = 1, 2, 3, 4, 5 indexes the most frequently occurring examples in the dataset. We have performed our simulations with these 25 different initial viral strains having the following point mutations: $G_{11} =$ $\{69D\}, G_{12} = \{70E\}, G_{13} = \{74I\}, G_{14} = \{151M\}, G_{15} = \{41L\}, G_{21} = \{69D, 115F\}, G_{22} = \{69D, 115F\}, G_{23} = \{69D, 115F\}, G_{24} = \{69D, 115F\}, G_{25} = \{60D, 115F\}, G_{25} = \{60D,$ $\{69D, 215Y\}, G_{23} = \{70R, 215Y\}, G_{24} = \{67N, 69D\}, G_{25} = \{67N, 70R\}, G_{31} = \{67N, 70R\}, G_{31} = \{67N, 70R\}, G_{32} = \{67N, 70R\}, G_{33} = \{67N,$ $\{69D, 115F, 215Y\}, G_{32} = \{69D, 70R, 115F\}, G_{33} = \{67N, 69D, 215Y\}, G_{34} = \{67N, 69D, 215Y\}, G_{35} = \{67N, 69D, 215Y\}, G_{35} = \{67N, 60D, 215Y\},$ $\{67N, 70R, 115F, 215Y\}, G_{43} = \{69D, 70R, 115F, 215Y\}, G_{44} = \{67N, 69D, 70R, 115F\},$ $G_{45} = \{65N, 69D, 70R, 215Y\}, G_{51} = \{65N\}, \{69D, 70R, 115F, 215Y\}, G_{52} = \{65N, 69D, 70R, 215Y\}, G_{51} = \{65N\}, \{69D, 70R, 215Y\}, G_{52} = \{65N, 69D, 70R, 215Y\}, G_{53} = \{65N, 69D, 70R, 215Y\}, G_{54} = \{65N, 69D, 70R, 215Y\}, G_{55} = \{65N, 60D, 70R, 215Y\}, G_{55} = \{65N, 70R, 215Y\}, G_{5$ $\{69D, 70R, 74F, 115F, 215Y\}, G_{53} = \{41L, 67N, 69D, 70R, 215Y\},\$ $G_{54} = \{65N, 67N, 69D, 70R, 215Y\}, G_{55} = \{67N, 69D, 70R, 74I, 215Y\}$. For instance, $G_{14} = \{151M\}$ strain has only one point mutation 151M and the rest of the amino acids are the same as wild type HIV-1 virus.

Measuring the therapy success

It is essential to track the success of the given antiretroviral therapy by hindering the viral dynamics from the AIDS phase, *i.e.*, by keeping the CD4 + T cell count as high as possible. The AIDS phase occurs when CD4 + T cell count is less than 200 cell/ μ l (*Kitahata et al., 2009*). Our primary criterion for the success of NRTI therapy is the occurrence and

nonoccurrence of the AIDS phase after initiation of the therapy with some initiation timing τ and adherence level α , as was done in cohort studies (*van Sighem et al., 2003*). Note that it is also possible to increase the CD4 + T cell counts of patients during the AIDS phase by the initiation of the ARTs (*Shoko & Chikobvu, 2019*). However, here we are not analyzing what happens after the AIDS phase. Our primary goal is to determine how evolutionary dynamics under the NRTI therapy affect the occurrence of the AIDS phase.

All simulations start with one infected CD4 + T cell and one infected macrophage cell with one of the initial strains G_{ij} . The simulation final time t_f is considered 20 years, and therapy success/failure is determined according to the occurrence of the AIDS phase in 20 years. However, we note that the clinical goal of ART therapy is the full suppression of detectable viremia. In our simulations, total suppression of detectable viremia is equivalent to not developing AIDS after 20 years. However, the opposite is false: detectable (> 200 copies/ml) suppression misses low copies of violent mutants, eventually leading to the AIDS phase. Therefore, we consider the AIDS occurrence as our output. In the clinic, therapy is redesigned if complete suppression is not observed. However, our simulations never redesign the treatment to distinguish between successful/failed drug combinations.

We run our simulations for randomly scattered 512 (α , τ) \in [0,1] × [0,365] pairs for predetermined initial strain G_{ij} . The success rate (*SR*) of a therapy is measured as the number of (α , τ) pairs that lead to protection from the AIDS phase in all 512 (α , τ) pairs. In Fig. 4, we show some representative simulation results of the multi-strain within-host model (5), starting with the $G_{51} = \{65N, 69D, 70R, 115F, 215Y\}$ strain under various mono and dual NRTI therapies with randomly scattered (α , τ) pairs. For this simulation setup, nine out of 21 NRTI therapy protocols have considerable success in preventing the patient from the AIDS phase. The importance of adherence level (α) and initiation timing (τ) is evident from the figure for all cases. In some cases, such as the DDI-D4T combination shown in Fig. 4, the initiation timing considerably affects the success rates. Higher τ values yield therapy failure even at high adherence levels. As observed from the figure, the D4T-3TC combination yields the best *SR* value by performing well for late initiation with perfect adherence levels. For the current case, the success of the D4T-3TC combination is mainly due to the behaviour of the therapy in the higher initiation timing (τ) region.

While the importance of the adherence levels is evident from its direct relation with infection rates, the importance of the initiation timing is non-evident and should be explained here clearly. In Fig. 5, we illustrate the effect of initiation timing τ in our multi-strain model (5) when initial strain and adherence level are selected as $G_{51} = \{65N, 69D, 70R, 115F, 215Y\}$ and $\alpha = 0.5$. According to Figs. 5A–5B, $\tau = 50$ yields successful therapy by maintaining healthy CD4 +T cell and macrophage cells at normal levels and declining the viral load to undetectable levels. On the other hand, when we assume the initiation timing as $\tau = 360$, virologic failure and AIDS phase are observed in Figs. 5C–5D. According to our model (5), the main difference between early and late initiation timing is the diversity of viral strains at the initiation to therapy times. Late initiation to the therapy increases the probability of the occurrence of the more resistant strains, even if their ancestors are slowly growing. For example, as we compare Fig. 5B with



Figure 4 Illustration of possible mono and dual NRTI therapy outcomes carried out using 512 random (α , τ) pairs in the current multi-strain within-host model (5). The initial strain has been selected as $G_{51} = \{65N, 69D, 70R, 115F, 215Y\}$. Blue circles represent the failure after 20 years of simulation, *i.e.*, the AIDS phase occurs when the patients start the therapy τ after infection and take the therapy with an adherence rate α . Purple squares mean that the therapy succeeds under the conditions mentioned above. *SR* values represent the success rate defined as SR = # of purple squares/# of all data points. Full-size \cong DOI: 10.7717/peerj.15033/fig-4

Fig. 5D, the two generations of mutant strains occur when $\tau = 360$ (Fig. 5D) while there exists only one generation of mutant strains when $\tau = 50$ (Fig. 5B). The two generations of mutant strains yield viral rebound and failure of the therapy in Fig. 5D.

If we go back to Fig. 4, the NRTI combinations having boundary lines with relatively low slope values are more sensitive to increasing values of τ since these therapies yield high variance in IC_{50} values of possible viral strains mutated from the initial strain. Therefore, in our modelling framework, the late initiation is directly related to the variance of IC_{50} values corresponding to the initial strain and possible mutants. Thus, the level and type of the NRTI therapy should be planned so that the reoccurrence of the viral strains should be blocked depending on the initiation time τ . Additionally, in the reoccurrence phase of viral strains, non-perfect adherence to the therapy leads to the selection of resistant strains (Fig. 5D). In this case, two possible problems arise:

- 1. If the therapy protocol of the patient is updated, therapy is less likely to be successful than when therapy was first started.
- 2. The probability of infecting another person with more resistant strains increases, and the probability of having an AIDS phase increases for the infected person.

The existence of low viral loads of new mutated strains is enough for selecting these strains after antiretroviral therapy. Therefore, according to our simulations, initiation timing is as crucial as the adherence level to overcome the AIDS phase and to protect the possible susceptible persons from more dangerous scenarios.



Figure 5 The effect of initiation timing is illustrated with healthy cell and virion counts. The initial strain is taken as $G_{51} = \{65N, 69D, 70R, 115F, 215Y\}$ and the common adherence level $\alpha = 0.5$ is considered. (A) Dynamics of T(t) and M(t) when $\tau = 50$, (B) dynamics of viral strains when $\tau = 50$, (C) dynamics of T(t) and M(t) when $\tau = 360$, (D) dynamics of viral strains when $\tau = 360$. Black dashed vertical lines in parts c and d denote the HIV detection limit in blood as 200 copies/ml (*Barletta, Edelman & Constantine, 2004*).

Full-size DOI: 10.7717/peerj.15033/fig-5

The NRTI mutants are known to have epistasis effects, which implies that the viral fitness of the mutant strain depends on the existing genetic background. The epistasis effects may lead to the selection of diverse branches in mutant generations (*Biswas et al., 2019*). Epistasis of mutations can impact the values of IC_{50} and fitness costs. The data we used to train our IC_{50} values implicitly includes epistatic effects. The *ANN* model that predicts IC_{50} values for mutants is expected to learn the epistatic interactions. However, it is not completely unlikely that some unobserved data may have unpredictable epistasis. Nevertheless, that variant being underrepresented in the data implies its irrelevance in the clinic. On the other hand, the fitness costs of mutants are assumed to be fixed due to a lack of enough data. Nevertheless, as we explain later, this assumption should not significantly impact our claims.

Simulation results

Here we have simulated our multi-strain within-host model (5) for all possible initial strains G_{ij} to observe the effect of initial strains on success rates. All possible mono and dual NRTI therapies have been implemented for randomly scattered 512 (α , τ) \in [0,1] × [0,365] pairs. The *SR* values of mono and dual NRTI therapies are calculated, and the well-performed combination results are comparatively illustrated in Fig. 6.

In line with Fig. 3 and Table 3, the D4T-3TC combination has been the best option for 20 out of 25 cases. The inhibitory potential of this combination is because of the pharmacokinetic parameters (see Table 1) of inhibitors, the drug-resistance profiles of



Figure 6 SR values of various NRTI combinations obtained by simulating multi-strain within-host model (5) with initial viral strain G_{ij} for randomly scattered 512 (α , τ) \in [0,1] \times [0,365] pairs. Full-size \cong DOI: 10.7717/peerj.15033/fig-6

inhibitors (see Table 3), and their Bliss-independent action on the target enzyme. Following the D4T-3TC combination, the TDF-D4T and D4T-AZT combinations are observed to be in first place in four and one out of 25 cases, respectively. The strong relation between the infection rate of an initial strain (and possible new strains) and the corresponding success rate value is evident from the correlation between Figs. 3 and 6. For instance, according to Fig. 3, the D4T-3TC combination yields fewer infection rates for most of the viral strains. Similarly, Fig. 6 shows that the D4T-3TC combination has great success rates for most of the initial viral strains. We will later quantitatively analyze the relationship between the infection rates of the detected viral strains and the success rates of the given therapies.

According to our modelling framework, since the fitness cost of all strains is assumed to be the same, the initial strain is dominant when the patient is diagnosed. Moreover, as evident from Figs. 5B–5D, considerable mutational variations at low copy numbers exist besides the initial strain. However, only the dominant strain is likely to be detected (strains having higher than 200 copies/ml in blood *Barletta, Edelman & Constantine, 2004*) when a phenosense assay is implemented. Thus, the clinician would only observe the initial strain and maybe a few mutational variations (according to Figs. 5B–5D, only the initial strain can be observed when the patient is diagnosed) to decide on the NRTI therapy protocol. Therefore, it is inevitable to ask whether the only predictor of the success rate is the detected viral strains at the diagnosis.

The undetected viral strains play a vital role in estimating the success rate and finding an optimal therapy protocol—especially their infection rates. We have trained regression models that predict therapy outcomes based on the infection rates of the initial strain and its mutants—the mutants will be referred to as first, second, third, fourth, and fifth generations. The first generation is mutated from the initial strain, whereas the second is mutated from the first. For the regression model, we aimed to determine how many



Figure 7 Prediction process of *SR* values from the infection rates of the detected and possible mutant strains. The models G_i are constructed by considering *i* generation of mutant strains and the detected strain itself. For each generation, mean and maximum values of the infection rates are assigned to the input of possible *ANN* and *MLR* models. *SR*_{ANN} and *SR*_{MLR} denote the *SR* prediction of the *ANN* and *MLR* models from the given infection rate input.

Full-size DOI: 10.7717/peerj.15033/fig-7

generations of the detected strain(s) should be considered to predict an optimal therapy. To answer this question quantitatively, we construct the *ANN* and *MLR* models for predicting the success rate of therapy from the infection rates of the existing mutant strains. We construct six *ANN* and *MLR* models denoted by G_i for $i = 0, 1, ..., 5.G_i$ denotes i - thgeneration of the detected strain(s) that has been considered in the inputs of the models. For instance, model G_0 only assumes the infection rates of the detected viral strain(s), and model G_3 considers the infection rates of the detected viral strain(s) and the first three-generation mutants of this strain(s). In each generation of mutant strains, we use two values: mean and maximum values of the infection rates of the considered generation. Thus, together with the detected viral strain, the model G_i has 2i + 1 dimensional input. 2iinput values denote the mean and maximum infection rates of i - th generation, and the remaining input value denotes the infection rate of the detected viral strain at the diagnosis. The graphical illustration of model G_i can be seen in Fig. 7.

Simulation results are given in Fig. 6 for 25 initial strains converted to the training data for the *ANN* and *MLR* models. 304 input–output relations have been obtained from various therapies having $SR \ge 0.02$. For the *ANN* models, this data is divided into the train, test, and validation sets (70%, 15%, and 15%). Each G_i model having the *ANN* architecture is trained using the scaled conjugate gradient algorithm. Similarly, for the *MLR* models, 20% of the data is considered as a test set, and the remaining 80% is used in the training process. To test the prediction performances of the *ANN* and *MLR* models, we have generated an external test dataset by simulating the model (5) with 25 random initial strains having one-to-five-point mutations, and 314 test sample is obtained. Additionally, to observe how well our *ANN* and *MLR* models classify the therapies as successful ($SR \ge 0.5$), the area under the receiving operating curves is measured for both the *ANN* and *MLR* models.

We illustrate the regression and classification performances of the *ANN* models on the training and test sets in Fig. 8. Figure 9 shows similar predictive performance metrics of the *MLR* models on both the training and test sets. The mean square error (*MSE*), linear



Figure 8 Regression and classification performances of models G_i having the ANN architectures on predicting the SR values of the therapies. Models G_i assume the infection rates of the detected strain and its first *i* mutant generations and have 2i + 1 input values. Mean square error (*MSE*), linear correlation coefficient (*R*), and area under the curve (*AUC*) metrics are presented for both training and test data. Full-size \square DOI: 10.7717/peerj.15033/fig-8

correlation coefficient (R), and area under the curve (AUC) metrics are presented for six G_i models having the ANN and MLR architectures. According to the test set performance of the models, model G_2 gives better MSE, R, and AUC values with both the ANN and MLR architectures. That means considering the infection rates of both the detected strains and the first two mutant generations of the detected strains led to better predictions.

On the other hand, the G_0 type models yield relatively poor regression and classification performances, *i.e.*, considering only the infection rate of the detected strains is not enough to estimate better therapy protocols. This implies that the possible undetected mutant generations should also be taken into account in determining the therapy protocols. Nevertheless, there is a threshold on the number of mutant generations that must be considered. Figure 8 (for ANN architectures) and Fig. 9 (for MLR architectures) show that models G_3 , G_4 and G_5 overfit the data and yield less accurate predictions than the model G_2 for both architectures. Additionally, for each G_i model, the ANN architecture yields a better approximation for the SR values than the MLR architectures.

DISCUSSIONS AND CONCLUSIONS

In this study, we have proposed a multi-strain within-host model of HIV infection with time-dependent NRTI therapy. Drug-resistant strains have been assumed to initiate the infection for the patients, and six available NRTI inhibitors with mono and dual combinations have been implemented in the simulations for various initiation timing and adherence levels. To assess the drug response curves with the IC_{50} values of the NRTI-resistant strains, artificial neural network models are trained for each inhibitor





by using the Stanford HIV drug resistance database. To describe time-drug efficiency and time-infection rate curves, pharmacokinetic parameters of the inhibitors have been calculated and hybridized with the corresponding IC_{50} values. We have designed our simulation environment to determine the effect of initial strains, initiation timing for the therapy protocol, and adherence levels to the given drug usage schedule on the occurrence of the AIDS phase within 20 years after infection.

According to our modelling framework, the success rate of the NRTI therapies in case of late initiation has led to the availability of more resistant viral strains, and then the resistant strains become dominant in the host plasma after an initial decline of the detected strain. Although some mathematical models assume implicitly that the initiation timing does not affect the success-failure of the therapy (*Dixit & Perelson, 2004; Rong, Feng & Perelson, 2007*), our multi-strain model catches the penalty of late initiation since the late initiation was proven to block the therapy success in various experimental results (*Kitahata et al., 2009; van Sighem et al., 2003*). Our simulation results have shown that in the case of the late initiation to therapy, the efficiency of the therapy should be far more than the early initiation case to prevent the possible AIDS phase.

We have shown that D4T-3TC, D4T-AZT, and TDF-D4T combinations are less likely to result in treatment failure. These inhibitors have been seen to provide fewer infection rates due to their pharmacokinetic parameters and IC_{50} values in the presence of various viral strains. According to our results, the success rate of accurately predicting the best therapy depended on the composition of detected strains and their possible further mutants. This observation implies that the emergence of new mutants from the initial strain is likely to have a considerable effect on the success of the therapy. Thus, it is more reasonable to

suggest the optimal therapy combinations for the patients by considering the detected viral strain and the undetected mutant, which most likely were generated from the detected strain.

The most important message of this article is that the undetected viral strains, at the diagnosis, may have considerable effects on therapy outcomes. Specifically, double mutants of the detected viral strain should be taken into account even if they were not detected. Earlier studies, such as Stanford HIVdb (*Talbot et al., 2010*), HIV-grade (*Obermeier et al., 2012*), REGA (*Van Laethem et al., 2002*), and ANRS (Agence Nationale de Recherches sur le SIDA, *Meynard et al. (2002)*) predicted the best possible therapy protocol. REGA is a rule-based model and was developed by scientists at Rega Institute for Medical Research and University Hospitals, and classifies the isolates as susceptible, intermediate, and resistant (*Van Laethem et al., 2002*). ANRS *Meynard et al. (2002)*; *Singh (2017)* is also a rule-based computational resistance classifier based on a linear combination of mutations. However, the undetected viral strains may lower the prediction power of such models. We have shown that a multi-strain within-host model (5) can help estimate undetected mutant strains and their role in optimal therapy selection.

A possible criticism of our model is that each mutant strain should have a unique fitness cost. However, we assume a constant factor for all mutants. To our best knowledge, there is not much data for specific strains to construct a machine-learning model as we did for the IC_{50} values. According to the theory, fitness costs can play a role in selecting resistant strains, which can alter our success rate. However, the fitness costs would affect the dynamics more at low drug concentrations. Luckily, the phase changes (AIDS or no AIDS) occur at relatively high adherence levels, which implies a relatively high concentration.

Our modeled treatments include only NRTIs, but current clinical practice includes additional drugs (*Aguilar et al., 2022*). Indeed, including the other components of ART would add to the realism. However, it is known that different classes of HIV drugs generally interact independently (*Rosenbloom et al., 2012; Jilek et al., 2012*). By the independence assumption, the relative ranking of NRTI therapies is relevant to consideration for ART. However, we would like to openly indicate that our model is not designed to suggest a better first line of treatment but rather to relatively rank NRTI combinations in a multiscale model.

This study has investigated the effect of NRTI inhibitors, which are the most important members of Highly Active Antiretroviral Therapy (HAART) (*Achhra & Boyd*, 2013). Since the Stanford drug resistance database also includes the genotype-phenotype data of protease inhibitors (PI), non-nucleotide reverse transcriptase inhibitors (NNRTI), and integrase inhibitors (II), some future studies may include these groups of inhibitors with possible mono, dual or triple drug combinations. Some existing HAART protocols may also be simulated through such a modelling framework. On the other hand, we have not considered the too-late initiation of the NRTI therapy at considerably low CD4 + T cell levels because of the failure of simulated therapy protocols in such situations. Some future works may

also investigate more comprehensive therapies to prevent patients from the AIDS phase when they are diagnosed too late.

ADDITIONAL INFORMATION AND DECLARATIONS

Funding

This work was supported by TUBITAK, 2232—International Fellowship for Outstanding Researchers, Project number 118C244. All the results are the sole responsibility of the authors. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Grant Disclosures

The following grant information was disclosed by the authors: TUBITAK, 2232—International Fellowship for Outstanding Researchers, Project number 118C244.

Competing Interests

The authors declare there are no competing interests.

Author Contributions

- Huseyin Tunc conceived and designed the experiments, performed the experiments, analyzed the data, prepared figures and/or tables, authored or reviewed drafts of the article, conceptualization; Data curation; Formal analysis; Investigation; Methodology; Resources; Software; Validation; Visualization; Writing original draft; Writing review & editing, and approved the final draft.
- Murat Sari conceived and designed the experiments, authored or reviewed drafts of the article, writing review & editing, and approved the final draft.
- Seyfullah Kotil conceived and designed the experiments, prepared figures and/or tables, authored or reviewed drafts of the article, conceptualization; Methodology; Visualization; Supervision, Writing review & editing, and approved the final draft.

Data Availability

The following information was supplied regarding data availability:

All data and necessary codes are available at Github and Zenodo: https://github.com/ tnchsyn/multistrainhivmodel; tnchsyn. (2023). tnchsyn/multistrainhivmodel: v1.0 (v1.0). Zenodo. https://doi.org/10.5281/zenodo.7547299.

Supplemental Information

Supplemental information for this article can be found online at http://dx.doi.org/10.7717/ peerj.15033#supplemental-information.

REFERENCES

Achhra A, Boyd M. 2013. Antiretroviral regimens sparing agents from the nucleoside(tide) reverse transcriptase inhibitor class: a review of the recent literature. *AIDS Research and Therapy* **10**:33.

- Aguilar G, Truong H, Ovelar P, Samudio T, Lopez G, Garca-Morales C, Tapia-Trejo D, Lpez-Snchez D, vila Ros S, Giron A, De Arias A, Rios-Gonzalez C, McFarland W.
 2022. HIV drug resistance in persons initiating or reinitiating first-line antiretroviral therapy in Paraguay: results of a national patient survey. *Journal of Medical Virology* 94(10):5061–5065 DOI 10.1002/jmv.27933.
- Alexaki A, Liu Y, Wigdahl B. 2008. Cellular reservoirs of HIV-1 and their role in viral persistence. *Current HIV Research* 6:388–400 DOI 10.2174/157016208785861195.
- Alshorman A, Al-hosainat N, Jackson T. 2022. Analysis of HIV latent infection model with multiple infection stages and different drug classes. *Journal of Biological Dynamics* 16(1):713–732 DOI 10.1080/17513758.2022.2113828.
- Amamuddy O, Bishop N, Bishop O. 2017. Improving fold resistance prediction of HIV-1 against protease and reverse transcriptase inhibitors using artificial neural networks. *BMC Bioinformatics* 18:369–376 DOI 10.1186/s12859-017-1782-x.
- **Ball C, Gilchrist M, Coombs D. 2007.** Modeling within-host evolution of HIV: mutation, competition and strain replacement. *Bulletin of Mathematical Biology* **69**:2361–2385.
- Barletta J, Edelman D, Constantine N. 2004. Lowering the detection limits of HIV-1 viral load using real-time immuno-PCR for HIV-1 p24 antigen. *American Journal of Clinical Pathology* 122:20–27.
- Biswas A, Haldane A, Arnold E, Levy R. 2019. Epistasis and entrenchment of drug resistance in HIV-1 subtype B. *Elife* 8:e50524 DOI 10.7554/eLife.50524.
- **Blower S, Aschenbach A, Gershengorn H, Kahn J. 2001.** Predicting the unpredictable: transmission of drug-resistant HIV. *Nature Medicine* **7**:1016–1020.
- **Chen W, Teng Z, Zhang L. 2021.** Global dynamics for a drug-sensitive and drug-resistant mixed strains of HIV infection model with saturated incidence and distributed delays. *Applied Mathematics and Computation* **406**:126284 DOI 10.1016/j.amc.2021.126284.
- Chun T, Davey R, Ostrowski M, Shawn Justement J, Engel D, Mullins JI, Fauci AS. 2000. Relationship between pre-existing viral reservoirsand the re-emergence of plasma viremia after discontinuation of highlyactive anti-retroviral therapy. *Nature Medicine* 6:757–761.
- Crowe S. 1995. Role of macrophages in the pathogenesis of human immunodeficiency virus (HIV) infection. *Australian and New Zealand Journal of Medicine* 25(6):777–783 DOI 10.1111/j.1445-5994.1995.tb02881.x.
- **Dixit N, Perelson A. 2004.** Complex patterns of viral load decay under antiretroviral therapy: influence of pharmacokinetics and intracellular delay. *Journal of Theoretical Biology* **226**:95–109.
- Doekes H, Fraser C, Lythgoe K. 2017. Effect of the latent reservoir on the evolution of HIV at the within- and between-host levels. *PLOS Computational Biology* 13(1):e1005228.
- Feng M, Sachs N, Xu M, Grobler J, Blair W, Hazuda D, Miller M, Lai M. 2016. Doravirine suppresses common nonnucleoside reverse transcriptase inhibitorassociated mutants at clinically relevant concentrations. *Antimicrobial Agents and Chemotherapy* 60(4):2241–2247.

- Haase A. 1999. Population biology of HIV-1 infection: viral and CD4+ T cell demographics and dynamics in lymphatic tissues. *Annual Review of Immunology* 17:625–656.
- Hadjiandreou M, Conejeros R, Vassiliadis V. 2007. Towards a long-term model construction for the dynamic simulation of HIV infection. *Mathematical Biosciences and Engineering* **4**(3):489–504.
- Hadjiandreou M, Conejeros R, Wilson D. 2009. Long-term HIV dynamics subject to continuous therapy and structured treatment interruptions. *Chemical Engineering Science* 64:1600–1617.
- Hendricks C, Cordeiro T, Gomes A, Stevenson M. 2021. The interplay of HIV-1 and macrophages in viral persistence. *Frontiers in Microbiology* 12:646447 DOI 10.3389/fmicb.2021.646447.
- Herbein G, Varin A. 2010. The macrophage in HIV-1 infection: from activation to deactivation? *Retrovirology* 7:33.
- Hernandez-Vargas E, Middleton R. 2013. Modeling the three stages in HIV infection. *Journal of Theoretical Biology* **320**(7):33–40.
- Hernandez-Vargas EA. 2019. *Modeling and control of infectious diseases: with MATLAB and R.* Berlin: Elsevier Academic Press.
- Holec A, Mandal S, Prathipati P, Destache C. 2018. Nucleotide reverse transcriptase inhibitors: a thorough review, present status and future perspective as HIV therapeutics. *Current HIV Research* 15(6):411–421.
- Jilek B, Zarr M, Sampah M, Rabi S, Bullen C, Lai J, Shen L. 2012. A quantitative basis for antiretroviral therapy for HIV-1 infection. *Nature Medicine* 18(3):456–465 DOI 10.1038/nm.2665.
- Kitahata M, Gange S, Abraham A, Merriman B, Saag M, Justice A, Hogg R, Deeks S, Eron J, Brooks J, Rourke S, Gill M, Bosch R, Martin J, Klein M, Jacobson L, Rodriguez B, Sterling T, Kirk G, Napravnik S, Rachlis A, Calzavara L, Horberg M, Silverberg M, Gebo K, Goedert J, Benson C, Collier A, Van Rompaey S, Crane H, McKaig R, Lau B, Freeman A, Moore R. 2009. Effect of early versus deferred antiretroviral therapy for HIV on survival. *AIDS* 360(18):1815–1826.
- Kruize Z, Kootstra N. 2019. The role of macrophages in HIV-1 persistence and pathogenesis. *Frontiers in Microbiology* 10:2828 DOI 10.3389/fmicb.2019.02828.
- Kuritzkes D. 2011. Drug resistance in HIV-1. *Current Opinion in Virology* 1(6):582–589 DOI 10.1016/j.coviro.2011.10.020.
- Khnert D, Kouyos R, Shirreff G, Peerska J, Scherrer AU, Bni J, Yerly S, Klimkait T, Aubert V, Gnthard HF, Stadler T, Bonhoeffer S, Study SHC. 2018. Quantifying the fitness cost of HIV-1 drug resistance mutations through phylodynamics. *PLOS Pathogens* 14:e1006895 DOI 10.1371/journal.ppat.1006895.
- Lagunin A, Rudik A, Pogodin P, Savosina P, Tarasova O, Dmitriev A, Ivanov S, Biziukova N, Druzhilovskiy D, Filimonov D, Poroikov V. 2023. CLC-Pred 2.0: a freely available web application for in silico prediction of human cell line cytotoxicity and molecular mechanisms of action for druglike compounds. *International Journal of Molecular Sciences* 24(2):1689 DOI 10.3390/ijms24021689.

- Lythgoe K, Pellis L, Fraser C. 2013. Is HIV short-sighted? Insights from a multi-strain nested model. *Evolution* 67(10):2769–2782.
- Masso M, Vaisman I. 2013. Sequence and structure based models of HIV-1 protease and reverse transcriptase drug resistance. *BMC Genomics* 14:S3.
- **Mauskopf J. 2013.** A methodological review of models used to estimate the cost effectiveness of antiretroviral regimens for the treatment of HIV infection. *Pharmacoeconomics* **31(11)**:1031–1050.
- Meynard J, Vray M, Morand-Joubert L, Race E, Descamps D, Peytavin G, Matheron S, Lamotte C, Guiramand S, Costagliola D, Brun-Vezinet F, Clavel F, Girard P. 2002. Phenotypic or genotypic resistance testing for choosing antiretroviral therapy after treatment failure: a randomized trial. *AIDS* 16(5):727–736.
- Montessori V, Press N, Harris M, Akagi L, Montaner J. 2004. Adverse effects of antiretroviral therapy for HIV infection. *CMAJ* 170(2):229–238.
- Obermeier M, Pironti A, Berg T, Braun P, Daumer M, Eberle J, Ehret R, Kaiser R, Kleinkauf N, Korn K, Kcherer C, Mller H, Noah C, Strmer M, Thielen A, Wolf E, Walter H. 2012. HIVGRADE: a publicly available, rules-based drug resistance interpretation algorithm integrating bioinformatic knowledge. *Intervirology* 55:102–107 DOI 10.1159/000331999.

Orenstein J. 2001. The macrophage in HIV infection. Immunobiology 204(5):598-602.

Perelson A, Nelson P. 1999. Mathematical analysis of HIV-1 dynamics in vivo. *SIAM Review* **41**(1):3–44.

- Pham H, Labrie L, Wijting I, Hassounah S, Lok K, Portna I, Goring M, Han Y, Lungu C, van der Ende M, Brenner B, Boucher C, Rijnders B, van Kampen J, Mesplde T, Wainberg M. 2018. The S230R integrase substitution associated with viral rebound during DTG monotherapy confers low levels INSTI drug resistance. *The Journal of Infectious Diseases* 218(5):698–706.
- Rhee S, Fessel W, Zolopa A, Hurley L, Liu T, Taylor J, Nguyen D, Slome S, Klein D, Horberg M, Flamm J, Follansbee S, Schapiro J, Shafer R. 2005. Protease and reverse-transcriptase mutations: correlations with antiretroviral therapy in subtype B isolates and implications for drug-resistance surveillance. *The Journal of Infectious Diseases* 192(3):456–465 DOI 10.1086/431601.
- Rhee S, Gonzales M, Kantor R, Betts B, Ravela J, Shafer R. 2003. Human immunodeficiency virus reverse transcriptase and protease sequence database. *Nucleic Acids Research* 31:298–303.
- Rhee S, Taylor J, Fessel W. 2010. HIV-1 protease mutations and protease inhibitor crossresistance. *Antimicrobial Agents and Chemotherapy* 54(10):4253–4261.
- Rong L, Feng Z, Perelson A. 2007. Emergence of HIV-1 drug resistance during antiretroviral treatment. *Bulletin of Mathematical Biology* **69**:2027–2060.
- Rosenbloom DI, Hill AL, Rabi SA, Siliciano RF, Nowak MA. 2012. Antiretroviral dynamics determines HIV evolution and predicts therapy outcome. *Nature Medicine* 18:1378–1385.
- **Rouzine I. 2022.** A role for CD4+ helper cells in HIV control and progression. *AIDS* **36(11)**:1501–1510 DOI 10.1097/QAD.00000000003296.

- Shah D, Freas C, Weber I, Harrison R. 2020. Evolution of drug resistance in HIV protease. *BMC Bioinformatics* 21:497–512 DOI 10.1186/s12859-020-03825-7.
- Shoko C, Chikobvu D. 2019. A superiority of viral load over CD4 cell count when predicting mortality in HIV patients on therapy. *BMC Infectious Diseases* 19(1):169 DOI 10.1186/s12879-019-3781-1.
- Van Sighem A, Van de Wiel M, Ghani A, Jambroes M, Reiss P, Gyssens I, Brinkman K, Lange J, De Wolf F. 2003. Mortality and progression to AIDS after starting highly active antiretroviral therapy. *AIDS* 17(15):2227–2236 DOI 10.1097/00002030-200310170-00011.
- Singh Y. 2017. Machine learning to improve the effectiveness of ANRS in predicting HIV drug resistance. *Healthcare Informatics Research* 23(4):271–276 DOI 10.4258/hir.2017.23.4.271.
- Steiner M, Gibson K, Crandall K. 2020. Drug resistance prediction using deep learning techniques on HIV-1 sequence data. *Viruses* 12(5):560 DOI 10.3390/v12050560.
- Sutimin S, Chirove F, Soewono E, Nuraini N, Suromo L. 2017. A model incorporating combined RTIs and PIs therapy during early HIV-1 infection. *Mathematical Biosciences* 285:102–111.
- Talbot A, Grant P, Taylor J, Baril J, Liu T, Charest H, Brenner B, Roger M, Shafer R, Cantin R, Zolopa A. 2010. Predicting tipranavir and darunavir resistance using genotypic, phenotypic, and virtual phenotypic resistance patterns: an independent cohort analysis of clinical isolates highly resistant to all other protease inhibitors. *Antimicrobial Agents and Chemotherapy* 54:2473–2479 DOI 10.1128/AAC.00096-10.
- Tarasova O, Biziukova N, Filimonov D, Poroikov V. 2018. A computational approach for the prediction of hiv resistance based on amino acid and nucleotide descriptors. *Molecules* 23(11):2751 DOI 10.3390/molecules23112751.
- Tarasova O, Biziukova N, Shemshura A, Filimonov D, Kireev D, Pokrovskaya A, Poroikov V. 2023. Identification of molecular mechanisms involved in viral infection progression based on text mining: case study for HIV infection. *International Journal of Molecular Sciences* 24(2):1465 DOI 10.3390/ijms24021465.

Tressler R, Godfrey C. 2012. NRTI backbone in HIV treatment. Drugs 72:2051–2062.

- Vaidya N, Rong L. 2017. Modeling pharmacodynamics on HIV latent infection: choice of drugs is key to successful cure via early therapy. *SIAM Journal on Applied Mathematics* 77:1781–1804.
- Valcour V, Chalermchai T, Sailasuta N, Marovich M, Lerdlum S, Suttichom D, Suwanwela N, Jagodzinski L, Michael N, Spudich S, van Griensven F, de Souza M, Kim J, Ananworanich J. 2012. Central nervous system viral invasion and inflammation during acute HIV infection. *The Journal of Infectious Diseases* 206(2):275–282.
- Van Laethem K, De Luca A, Antinori A, Cingolani A, Perno C. 2002. A genotypic drug resistance interpretation algorithm that significantly predicts therapy response in HIV-1-infected patients. *Antiviral Therapy* 7:123–129.
- Wu P, Zhao H. 2020. Dynamics of an HIV infection model with two infection routes and evolutionary competition between two viral strains. *Applied Mathematical Modelling* 84:240–264.

Zhang J, Rhee S, Taylor J, Shafer R. 2005. Comparison of the precision and sensitivity of the antivirogram and PhenoSense HIV drug susceptibility assays. *JAIDS: Journal of Acquired Immune Deficiency Syndromes* **38**:439–444.