

## References

- [1] HIV Sequence Database (2015) Los Alamos National Laboratory. <http://www.hiv.lanl.gov/>. Accessed June 24, 2015.
- [2] Simmonds P, et al. (1991) Discontinuous sequence change of human immunodeficiency virus (HIV) type 1 env sequences in plasma viral and lymphocyte-associated proviral populations in vivo: implications for models of HIV pathogenesis. *Journal of Virology* 65(11):6266–6276.
- [3] Shankarappa R, et al. (1999) Consistent viral evolutionary changes associated with the progression of human immunodeficiency virus type 1 infection. *J. Virol.* 73(12):10489–10502.
- [4] Edwards CTT, et al. (2006) Population genetic estimation of the loss of genetic diversity during horizontal transmission of HIV-1. *BMC Evolutionary Biology* 6(28).
- [5] Fischer M, et al. (2004) Attenuated and nonproductive viral transcription in the lymphatic tissue of HIV-1-infected patients receiving potent antiretroviral therapy. *J. Infect. Dis.* 189(2):273–285.
- [6] Novitsky V, et al. (2009) Timing constraints of in vivo gag mutations during primary HIV-1 subtype C infection. *PLoS ONE* 4(11):e7727.
- [7] Akaike H (1974) A new look at the statistical model identification. *IEEE Transactions on Automatic Control* 19(6):716–723.

Table S1: Summary of the patient data collected from the LANL HIV sequence database [1] in the data sets from public sources.

Reference	Patient ID	Sequences		Time points		Time span		$\Delta$ AIC	Root	MAE		MAD	
		Plasma	PBMC	Total	Plasma	PBMC	Total	Plasma		Years	Scaled	Years	Scaled
[2]	2658	69		69	5	5	5	5.2	-1.2	0.42	0.081	0.64	0.12
[3]	825	49		49	6	6	6	8.2	0.11	0.73	0.089	0.89	0.11
[4]	<b>7259</b>	<b>28</b>		<b>28</b>	<b>4</b>	<b>4</b>	<b>4</b>	<b>2.4</b>	<b>-0.97</b>	<b>0.69</b>	<b>0.28</b>	<b>0.61</b>	<b>0.25</b>
	<b>7265</b>	<b>21</b>		<b>21</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>0.75</b>	<b>-0.73</b>	<b>0.18</b>	<b>0.25</b>	<b>0.33</b>	<b>0.44</b>
	13333	38		38	4	4	4	1.4	-0.49	0.31	0.22	0.30	0.22
	13334	36		36	5	5	5	2.0	-1.2	0.47	0.22	0.44	0.22
	13336	42		42	4	4	4	2.0	-0.31	0.34	0.17	0.34	0.17
[3]	<b>820*</b>	<b>45</b>	<b>81</b>	<b>126</b>	<b>5</b>	<b>10</b>	<b>15</b>	<b>2.1</b>	<b>3.3</b>	<b>0.78</b>	<b>0.38</b>	<b>8.2</b>	<b>3.9</b>
	821	69	178	247	7	17	17	6.5	-0.20	0.46	0.070	0.78	0.12
	822	29	90	119	3	10	10	5.8	-1.6	0.50	0.086	0.89	0.15
	824	52	102	154	7	9	13	8.6	0.87	0.69	0.080	1.3	0.15
	10138*	64	112	176	6	13	16	8.0	-1.0	1.5	0.19	1.2	0.15
	13889	77	65	142	13	14	18	13	-3.5	1.5	0.11	1.8	0.13
[5]	<b>10769</b>	<b>108</b>	<b>56</b>	<b>164</b>	<b>10</b>	<b>4</b>	<b>11</b>	<b>5.4</b>	<b>5.8</b>	<b>5.60</b>	<b>1.0</b>	<b>7.6</b>	<b>1.4</b>
[6]	34391	12	35	52	3	5	6	0.91	-0.99	0.12	0.14	0.24	0.27
	<b>34399*</b>	<b>50</b>	<b>72</b>	<b>122</b>	<b>5</b>	<b>9</b>	<b>12</b>	<b>2.5</b>	<b>0.38</b>	<b>0.31</b>	<b>0.12</b>	<b>0.63</b>	<b>0.25</b>
	34411	14	19	33	3	3	6	1.3	-0.027	0.14	0.11	0.25	0.20

Patient ID corresponds to the anonymized patient identifiers in the LANL database. Time span is in years.  $\Delta$ AIC is the Akaike Information Criterion (AIC) [7] of the null model minus the AIC of the linear model. Root is the estimate of the root time in years by the linear model with respect to the time of the first sample. MAE is the mean absolute error (between collection date and estimated date) of the training data. MAD is mean absolute difference (between collection date and estimated date) of the censored data. Scaled MAE(D) is the MAE(D) divided by the time span of the training data. The model fail to be calibrated for bold rows.