

# *sierra-local*: A lightweight standalone application for secure HIV-1 drug resistance prediction

Jasper C Ho<sup>a</sup>, Garway T Ng<sup>a</sup>, Mathias Renaud<sup>a</sup>, Art FY Poon<sup>a,b#</sup>

<sup>a</sup>Department of Pathology and Laboratory Medicine, Western University, London, ON, Canada;

<sup>b</sup>Department of Microbiology and Immunology, Western University, London, ON, Canada.

## 1 Abstract

2 Genotypic resistance interpretation systems for the prediction and interpretation of HIV-1  
3 antiretroviral resistance are an important part of the clinical management of HIV-1 infection.  
4 Current interpretation systems are generally hosted on remote web servers that enable clinical  
5 laboratories to generate resistance predictions easily and quickly from patient HIV-1 sequences  
6 encoding the primary targets of modern antiretroviral therapy. However they also potentially  
7 compromise a health provider's ethical, professional, and legal obligations to data security,  
8 patient information confidentiality, and data provenance. Furthermore, reliance on web-based  
9 algorithms makes the clinical management of HIV-1 dependent on a network connection. Here,  
10 we describe the development and validation of *sierra-local*, an open-source implementation of  
11 the Stanford HIVdb genotypic resistance interpretation system for local execution, which aims  
12 to resolve the ethical, legal, and infrastructure issues associated with remote computing. This  
13 package reproduces the HIV-1 resistance scoring by the web-based Stanford HIVdb algorithm  
14 with a high degree of concordance (99.997%) and a higher level of performance than current  
15 methods of accessing HIVdb programmatically.

## INTRODUCTION

16 Genotype-based prediction of human immunodeficiency virus type 1 (HIV-1) drug resistance is an  
17 important component for the routine clinical management of HIV-1 infection [1, 2]. Detecting the  
18 presence of viruses carrying mutations that confer drug resistance enables physicians to select an  
19 optimal drug combination for that patient's treatment regimen. Furthermore, genotyping by bulk

---

# Corresponding author, apoon42@uwo.ca

20 sequencing is a cost-effective alternative to the direct measurement of drug resistance from cultur-  
21 ing virus isolates in a laboratory [3]. Provided access to affordable bulk sequencing at an accredited  
22 laboratory for clinical microbiology, the interpretation of HIV-1 sequence variation is the primary  
23 obstacle to utilizing resistance genotyping for HIV-1 care. Fortunately, there are several HIV-1  
24 drug resistance interpretation algorithms that can be accessed at no cost through web applications  
25 or services hosted by remote network servers, such as the Standard University HIV Drug Resis-  
26 tance Database (HIVdb) [4], Agence Nationale de Recherche sur le SIDA (ANRS) AC11 [5], and  
27 Rega Institute [6] algorithms. The Stanford HIVdb interpretation system can be accessed either  
28 through a web browser at <http://hivdb.stanford.edu/hivdb> or programmatically through its Sierra  
29 Web Service [7], which requires the transmission of an HIV-1 sequence from a local computer  
30 over the network to the remote server. This is a convenient arrangement for clinical laboratories  
31 because there is no need to install any specialized software, web browsers are ubiquitous and most  
32 users are familiar with submitting web forms.

33 On the other hand, there are a number of disadvantages to accessing interpretation systems over  
34 a network connection. First, HIV-1 sequences are sensitive patient information, not only because  
35 infection with HIV-1 remains a highly stigmatized condition, but also because sequence data have  
36 been used as evidence in the criminal prosecution of individuals for engaging in sexual intercourse  
37 without disclosing their infection status, leading to virus transmission [8]. Once sequence data  
38 have been transmitted to a remote server, one cedes all control over data security. Preventing the  
39 onward distribution of the data and deleting the data once the analysis is complete, for instance,  
40 is entirely the responsibility of the system administrators of the host server. Furthermore, unless  
41 the host server employs a secure transfer protocol, the unencrypted data are transmitted in the  
42 clear between a number of intermediary web servers, exposing these data to a ‘man-in-the-middle’  
43 attack [9].

44 Second, the algorithm hosted on the server is effectively a black box — one has no insight  
45 into how resistance predictions are generated. Even if a version of the algorithm has been released  
46 into the public domain, one cannot be certain that the exact same algorithm was applied to their

47 transmitted data. Importantly, different versions of a given algorithm can output significantly dif-  
48 ferent resistance predictions, with the general trend being an increase in both resistance scores and  
49 predicted resistance levels [10]. In addition to contributing to inconsistencies in algorithm outputs,  
50 this makes it difficult to track data provenance, *i.e.*, the historical record of data processing, that  
51 has become recognized as a critical gap in the workflows of clinical laboratories. For instance, the  
52 College of American Pathologists recently issued new accreditation requirements stipulating that  
53 clinical laboratories must track the specific version of software programs used to process patient  
54 data [11]. Thus, a reliance on web-based systems creates significant issues for the reproducibility  
55 and quality assurance of clinical workflows. The Stanford HIVdb web service (Sierra [7]), for  
56 instance, automatically utilizes the most recent version of the HIVdb algorithm. While this con-  
57 straint ensures that users employ the most up-to-date algorithm, it also introduces hidden changes  
58 to clinical pipelines, which may have been locally validated on older versions of the algorithm.

59 Third, dependence on a web resource may cause problems when the laboratory cannot access  
60 the host server, either due to local or regional network outages, or because the host server is mal-  
61 functioning or offline. In our experience, the web servers hosting the more popular HIV drug  
62 resistance interpretation algorithms such as the Stanford HIVdb database are reliable and well-  
63 maintained. However, it is not unusual for other web-based algorithms to be relocated or go offline  
64 when the developers move to other institutions or lack the resources to maintain the service.

65 One of the important features of the Stanford HIVdb algorithm is that it is regularly updated and  
66 released into the public domain in a standardized XML-based interchange format — the Algorithm  
67 Specification Interface version 2 (ASI2) format [12] — that was formulated and published by  
68 the same developers in conjunction with the Frontier Science Foundation. Here, we describe the  
69 implementation and validation of *sierra-local*, an open-source Python package for local execution  
70 of the HIVdb algorithm in the ASI2 format. This package utilizes, but does not require, a network  
71 connection to synchronize its local ASI2 file and reference data with the latest releases on the  
72 Stanford HIVdb web server. Our objective was to release a lightweight alternative to transmitting  
73 HIV-1 sequences to the HIVdb web server that minimizes the number of software dependencies,

74 and that produces the exact same interpretations as the Sierra web service for all available HIV-1  
75 sequences in the Stanford database.

## MATERIALS AND METHODS

76 **Data Collection.** We obtained the entirety of the genotype-treatment correlation datasets available  
77 through the Stanford HIV Drug Resistance Database (HIVdb [13]) on May 7 2018. These data  
78 included nucleotide sequences of 105,694 protease (PR) isolates, 112,723 reverse-transcriptase  
79 (RT) isolates, and 12,332 integrase (IN) isolates for a total of 230,749 sequences. In addition  
80 to sequence data, each record comprised of a list of the specific antiretrovirals (ARVs) that each  
81 isolate had been exposed to *in vivo* prior to collection, the region and year of collection, and  
82 subtype as determined by the Stanford University HIV Drug Resistance Database's HIV Subtyping  
83 Program. After screening for empty and invalid data, the resulting dataset contained 103,711 PR  
84 entries, 110,222 RT entries, and 11,769 IN entries totalling 226,702 records. In addition, we  
85 retrieved 7 population-based HIV-1 *pol* datasets from Genbank using the NCBI PopSet interface  
86 (<http://www.ncbi.nlm.nih.gov/popset>). These datasets were selected from the most recent uploads  
87 of substantial numbers of HIV-1 sequences covering the regions encoding both PR and RT, and  
88 representing a diversity of HIV-1 subtypes and sampling locations around the world.

89 **Local HIVdb Algorithm Implementation.** The team at Stanford HIVdb have created a Python  
90 tool, SierraPy (<https://github.com/hivdb/sierra-client/tree/master/python>), that serves as a comm-  
91 and-line interface (CLI) for HIVdb. SierraPy does not process sequences directly, however, and  
92 only serves as a front-end for the HIVdb Sierra GraphQL Web Service (<https://hivdb.stanford.edu/graphql>) [7]. Its reliance on an active network connection to offload sequence processing to  
93 a remote server does not fulfill the usage gap we aim to address with *sierra-local*. To fill this  
94 gap, a system that provides a complete interface to a local version of the algorithm is needed.  
95 This local algorithm is first obtained as a publicly available HIVdb ASI2 file, which encodes both  
96 the algorithm for resistance scoring sequences and annotations describing relevant drug resistance  
97 mutations (DRMs) and ARVs of interest [12]. In short, it serves as a container for the core of the

99 HIVdb resistance interpretation system which is not directly usable in a data-processing pipeline,  
100 as it is essentially only a descriptive XML defining the rules by which sequences should be scored.  
101 In the existing HIVdb pipeline, a Java-based interpreter generator called SableCC is used to com-  
102 pile an algorithm interpreter from the HIVdb ASI2 file, but we have not been able to find any such  
103 compiler in Python. The usage of SableCC in our local implementation of the Stanford algorithm  
104 would introduce further dependencies and obfuscate the clear relationship between the algorithm  
105 file and how local interpretations are generated. Hence, in light of the need for a Python-based  
106 local interpreter for the ASI2 format, we developed a regular expressions-based keyword-parsing  
107 method by which *sierra-local* locally compiles an executable model in Python directly from the  
108 local algorithm file. This method iterates through the HIVdb algorithm as an XML tree object in  
109 Python 3 and extracts the information encoded within using ASI2 keywords defined by the ASI2  
110 Document Type Definition (DTD). *sierra-local* then uses this method to calibrate the model by  
111 assigning the drug clause definitions, drug class lists, resistance level interpretations, DRM com-  
112 ments, and complex drug-DRM scoring conditions to a set of dictionary and list objects. Once  
113 populated, this model serves as the framework for sequence resistance scoring.

114 **Sequence Pre-Processing and Validation.** Prior to scoring, the HIVdb Sierra Web Service per-  
115 forms several pre-processing and validation steps on submitted query sequences and identified  
116 mutation sites found within sequences. We emulated these steps to maximize fidelity to the HIVdb  
117 pipeline, including sequence alignment, gene identification, mutation site classification, sequence  
118 trimming, sequence subtyping, and sequence validation. These pre-processing steps can be consid-  
119 ered parts of the algorithm involved in generating resistance prediction scores that are not included  
120 and distributed in the HIVdb ASI2 file.

121 **Sequence Alignment.** Of particular importance in these steps is sequence alignment. HIVdb uti-  
122 lizes NucAmino [14] to initially align and identify amino acid mutations in each query sequence.  
123 This nucleotide-to-amino acid alignment program is optimized for viral gene sequences, and was  
124 developed by the HIVdb developers in the Go language. However, in practice, NucAmino does  
125 not return the aligned sequences themselves; instead it only returns a list of mutations relative to

126 the consensus subtype B amino acid sequence and general sequence metadata, *e.g.*, the aligned  
127 start and end coordinates of the query sequence relative to HIV-1 *pol*. If *sierra-local* does not  
128 align the query viral sequences exactly as how HIVdb would, the mutations identified relative to  
129 the consensus subtype B sequence would not be identical in all cases. Since these mutations serve  
130 as the basis for resistance predictions, an identical alignment process is of the utmost importance  
131 in maintaining fidelity to HIVdb Sierra. Thus, we decided to incorporate the NucAmino align-  
132 ment program as a dependency, rather than substituting a more integrated native implementation in  
133 Python. *sierra-local* calls a pre-compiled NucAmino binary as a Python subprocess with default  
134 settings to execute this pre-processing step. We used NucAmino version 0.1.3 for our validation  
135 experiments. The optional JSON output from NucAmino was captured in Python by redirecting  
136 the standard output stream from the subprocess.

137 *Gene Identification.* A critical part of resistance scoring is knowing which gene products encoded  
138 by HIV-1 *pol* (protease, RT and integrase) are present in the sequence being analyzed. To map  
139 the query sequence to these targets, we compared the aligned start and end positions returned by  
140 NucAmino to the HXB2 reference positions. For consistency with the HIVdb pipeline, amino  
141 acids were renumbered relative to the start position of the corresponding gene product (PR, RT  
142 and IN).

143 *Mutation Site Classification.* In the process of mutation site classification, each amino acid mu-  
144 tation site identified by NucAmino was further categorized as an insertion, deletion, or mixture  
145 using the original nucleotide sequence as a reference for amino acid translation. For consistency,  
146 we ported a Java method from the Sierra algorithm for determining nucleotide codon translation  
147 ambiguity to Python. Each site of interest identified by NucAmino on the aligned and translated  
148 query sequence generated a list of characters representing possible encoded amino acids at that site.  
149 Associating mutation sites with more than one encoded amino acid allows for a ‘fuzzy’ matching  
150 of a single sequence to multiple scoring conditions sharing the same residue position but different  
151 mutations. In certain cases, ‘highly ambiguous’ sites encoding more than four possible amino acids  
152 – made possible due to the presence of ambiguous nucleotides such as ‘R’ (A/G), ‘B’ (C/G/T), and

153 ‘N’ (A/T/C/G) – were flagged as ‘ambiguous’ and represented with a translation of ‘X’.

154 **Sequence Trimming.** Leading and trailing regions in each sequence containing a minimum pro-  
155 portion (30%) of ‘low quality’ sites — defined as sequenced sites that are highly ambiguous, stop  
156 codons, unusual mutations, or frameshifts — with at least one of these sites every 15 residues were  
157 trimmed prior to resistance scoring. Based on our inspection of the Sierra source code, we defined  
158 a site as sequenced if the codon does not have more than one unknown nucleotide. Sites further  
159 qualified as ‘unusual mutations’ if they were indicative of APOBEC-mediated G-to-A hypermu-  
160 tation [15] or if the highest frequency of that mutation in the pooled untreated and treated viruses  
161 for that specific group M subtype was less than 0.1% in the Stanford University HIV Drug Resis-  
162 tance Database. We configured the *sierra-local* installation process to obtain a local copy of the  
163 reference data for APOBEC-mediated G-to-A hypermutations in HIV-1 *pol* and for other HIV-1  
164 mutation prevalences from the Stanford HIVdb server, and to automatically update this local copy  
165 to accommodate changes in these reference data over time.

166 **Sequence Subtyping.** The previously discussed trimming step requires that sequences be subtyped  
167 in order to determine the frequency of mutations in subtype-specific pooled untreated and treated  
168 viruses. We wrote a Python implementation of HIVdb’s HIV Subtyping Program, which cate-  
169 gorizes submitted sequences as a pure subtype, a circulating recombinant form (CRF), a unique  
170 recombinant form (URF), non-group M HIV-1, or HIV-2. This process calculates uncorrected  
171 pairwise distances (the proportion of nucleotide differences) between submitted nucleotide se-  
172 quences and a set of 200 different subtype-specific reference sequences. Because the HIVdb HIV  
173 Subtyping Program is very simple and does not perform phylogenetic analysis or bootstrapping,  
174 its results may not be as accurate as more sophisticated systems more commonly used such as  
175 the Rega Institute HIV-1 Automated Subtyping Tool [4]. Comparing interpretation systems, it  
176 has been suggested that ANRS AC11, HIVdb, and Rega demonstrate discordances that may be  
177 subtype-dependent [16, 17]. Yet, current genotypic resistance interpretation systems, including  
178 Stanford HIVdb, are subtype-agnostic within themselves, meaning that they do not offer differen-  
179 tial resistance penalty scoring based on the subtype identified. Thus, other than being used to trim

180 submitted sequences, the subtyping in this system does not significantly influence HIV-1 sequence  
181 interpretation.

182 **Resistance Scoring.** The HIVdb algorithm begins the scoring process by assessing each se-  
183 quence's potential resistance to 22 commonly-used ARVs independently by searching for ARV-  
184 specific mutation conditions present in the query sequence. ARV-specific mutation conditions  
185 encode and define the circumstances in which a mutation in a particular gene has been shown to  
186 influence resistance, and quantifies this resistance with an integer penalty score. Each of these  
187 mutation conditions are associated with one or more specified amino acid changes at one or more  
188 particular positions, the ARV of interest, and the gene of interest in HIV-1 *pol* targeted by the ARV.  
189 For a mutation condition's resistance score to be counted, all of these criteria must be fulfilled.

190 *sierra-local* iterates over the ARV-specific mutation conditions corresponding to the gene re-  
191 gions detected in the query sequence. HIVdb Sierra also validates sequences and returns a list of  
192 validation problems found in each sequence. A sequence is invalidated or flagged as a query if: no  
193 genes are found; the sequence is a reverse complement; the genes have not been aligned properly;  
194 it is too short based on gene-specific minimum cutoff nucleotide lengths; it is trimmed at the 5' or  
195 3' ends due to low-quality leading or trailing sites; an indel is longer than 30 base pairs; invalid  
196 'NA' characters are found; one or more stop codons are found; too many 'unusual mutations' as  
197 previously described are found; the virus is subtyped as HIV-2; the number of APOBEC mutations  
198 is two or greater; the number of APOBEC mutations at drug resistance positions (DRPs) is one or  
199 greater; or if the count of frameshifts, unusual insertions, and unusual deletions together is positive.  
200 We emulated these validation steps in our pipeline to maximize parity with Sierra.

201 **Algorithm Output Generation.** Once all queries are pre-processed, scored, and otherwise anno-  
202 tated, *sierra-local* writes these results into a JSON (JavaScript Object Notation) format that mimics  
203 the standard output format of the HIVdb Sierra Web Service. For the sake of brevity and simplic-  
204 ity, we decided to have *sierra-local* omit the 'pretty pairwise' sequence output found in Sierra's  
205 standard output. This output format usually contains a numerical sequence of all residue positions,  
206 a reference amino acid sequence, an aligned nucleotide sequence, and a mutation-only list. These

207 data are not critical to a rapid genotypic interpretation system and may be omitted without detri-  
208 ment to the central aim. Furthermore, their exclusion from the results leads to a five-fold reduction  
209 in the results file size, which may be significant when processing large batches of sequences.

210 **Validation Against HIVdb Sierra.** We scored the entirety of the partitioned and filtered HIVdb  
211 genotype–treatment correlation dataset as of May 7 2018 with both *sierra-local* and SierraPy (ver-  
212 sion 0.2.1), storing all output results files from both programs. Validation was conducted using the  
213 HIVdb version 8.5 algorithm on both platforms. Because the algorithm was updated to version  
214 8.6.1 during the validation experiments, we used the newer version for the HIV-1 integrase data  
215 sets since the update mostly affected the interpretation of mutations within this region. Subsequent  
216 analysis of validation testing was conducted in R (version 3.4.4) [18] using the *jsonlite* package  
217 (version 1.5) [19] to extract resistance scores and relevant meta-data from the JSON results files  
218 from either program. For each sequence record, the resistance scores from the pipelines were com-  
219 pared for identicity. Validation and mutational metadata for discordant cases were analyzed for  
220 iterative refinement of the *sierra-local* scripts until a satisfactory level of concordance was attained.  
221 Subsequently, we carried out a second set of validation experiments on longer HIV-1 *pol* (PR-RT)  
222 sequence data sets under HIVdb algorithm version 8.7.

223 **Performance Measurements.** To quantify the speed of *sierra-local* in processing sequences, we  
224 chose a random subsample of gene-specific batches balanced across the genes of HIV-1 *pol* in  
225 our filtered genotype-treatment correlation dataset. This sampled dataset comprised of 9,673 PR  
226 sequences, 10,000 RT sequences, and 9,720 IN sequences contained in 10 batches per gene. Each  
227 gene-specific batch, roughly 10,000 sequences in size, was independently processed with both  
228 *sierra-local* and SierraPy. This method was timed using the *time* package in Python to determine  
229 file run-times and the resulting processing speeds of each program, measured in sequences per  
230 second.

231 **Software Availability.** The source code for *sierra-local* has been released under the GNU General  
232 Public License (version 3) and may be obtained at <http://www.github.com/PoonLab/sierra-local> or

233 from the Python Package Index (<https://pypi.org/project/sierralocal>). Detailed installation instruc-  
234 tions are provided on the GitHub website.

## RESULTS AND DISCUSSION

235 **Implementation.** Distribution of the Stanford HIVdb algorithm in an XML-based interchange  
236 format called the Algorithm Specification Interface version 2 (ASI2) enables a seamless approach  
237 to updating the algorithm and algorithm version control. Despite the accessibility of the HIVdb  
238 algorithm itself, the ASI2 file is necessary but not sufficient to generate resistance predictions. Ad-  
239 ditional steps needed but not encoded by the algorithm file include: sequence alignment, sequence  
240 quality control and validation, sequence trimming, sequence subtyping, and formatting of the re-  
241 sults output file with resistance predictions and accessory meta-data. These processes comprise  
242 the bulk of the functionality developed in *sierra-local* and are coordinated to generate resistance  
243 predictions in a manner as identical as possible to HIVdb Sierra.

244 As new versions of HIVdb are released to reflect the growing knowledge of HIV resistance,  
245 users may easily manually update their local copy of the algorithm with a provided Python script  
246 (`updater.py`). We also provide the option to automatically update the algorithm to the most recently  
247 released version at run-time. For example, three major updates to the HIVdb algorithm were  
248 released to the public during the development of *sierra-local*. Our local copies of these files were  
249 automatically retrieved by the *sierra-local* pipeline, but we also configured the pipeline to use  
250 specific versions by setting the ‘-xml’ option. The option to choose between automatic updating or  
251 freezing the algorithm to a specific version enables physicians and researchers to fulfill potential  
252 version tracking and data provenance requirements. However, even with functionalities addressing  
253 these obligations embedded in *sierra-local*, the health care provider still bears the responsibility  
254 of operating with a knowledge of their software and algorithm versions and the changes between  
255 these.

256 **Concordance with HIVdb Sierra.** Out of the 103,711 PR, the 111,222 RT, and the 11,769 IN  
257 sequences (total 226,702) processed with both *sierra-local* and HIVdb SierraPy pipelines, the pre-

258 dictated resistance scores and component subscores were completely identical in 226,696 (99.997%).  
259 Of the 6 sequences that did not have identical scores for all ARVs between the pipelines, 3 were  
260 PR sequences and 3 were RT (Table 1). The most frequent cause of discordance was the trimming  
261 of nucleotides in the leading or trailing ends of the sequence on the basis of the prevalence of the  
262 amino acid polymorphism in the corresponding HIV-1 subtype or ambiguous base calls.

263 We examined the distributions of sequence lengths and mixtures (ambiguous base calls) in  
264 the database to determine whether the discordant cases might be explained by these factors. Of  
265 the three discordant PR sequences, AY739171 contained an extremely large number of mixtures  
266 (15). In addition, this sequence was derived from an HIV-1 group O infection. Only 0.42% of  
267 PR sequences in the databases contained as many or more mixtures; the median [2.5% and 97.5%  
268 quantiles] number of mixtures was 1 [0, 9]. The remaining two PR sequences were unusually  
269 short, where the median length was 297 nt; the proportion of sequences with lengths shorter than  
270 GU188744 and JQ028402 were 0.43% and 0.76%, respectively. Although *sierra-local* reported  
271 non-zero resistance scores where *SierraPy* reported none, the scores were generally in the range  
272 of susceptible to potential low-level resistance interpretations. Similarly, the three discordant RT  
273 sequences were either short (overall median [quantiles] = 774 [588, 1680] nt) or contained substan-  
274 tial numbers of mixtures for sequences of comparable length (e.g., 4 [0, 19] mixtures for sequences  
275 between 550 and 650 nt in length). In the latter case, however, neither KT745612 nor KC221011  
276 contained a significantly excessive number of mixtures. These three RT sequences also resulted  
277 in slightly more discordant resistance interpretations; for example, the d4T resistance score for  
278 KC221011 was switched by *sierra-local* from intermediate to high-level resistance.

279 In addition, we ran both pipelines on six recently published sets of HIV-1 *pol* sequences com-  
280 prising both PR and RT encoding regions. These data sets were selected to cover a diversity of  
281 HIV-1 subtypes and locations around the world (Table 2). The major HIV-1 group M subtypes  
282 A, B, C and D were represented in these data, as well as several circulating recombinant forms  
283 (CRFs) such as CRF07\_BC, which is highly prevalent in East Asia. All resistance scores for all  
284 1,006 sequences were completely concordant between the pipelines.

Accession	Gene	Length (nt)	Mixtures	Reason	Drug	<i>SierraPy</i>	<i>sierra-local</i>
AY739171	PR	279	15	No genes found			
GU188744	PR	216	0	Sequence trimming, unusual mutations	ATV/r DRV/r FPV/r IDV/r LPV/r NFV SQV/r TPV/r	0 0 0 0 0 0 0 0	5 5 10 5 5 10 5 10
JQ028402	PR	243	0	Sequence trimming	ATV/r FPV/r IDV/r NFV SQV/r	0 0 0 0 0	5 5 5 15 5
KT745612	RT	630	10	Stop codon	ABC AZT d4T DDI TDF	100 130 130 100 70	95 125 125 95 65
DQ297313	RT	564	0	Sequence trimming	ABC AZT d4T DDI TDF	0 0 0 0 0	5 20 20 10 5
KC221011	RT	597	14	Sequence trimming	ABC AZT d4T DDI TDF	10 25 35 45 10	25 50 60 65 25

**Table 1** Discordant cases between *SierraPy* and *sierra-local*. Both pipelines were applied to the same database, comprising 103,711 PR, 110,222 RT and 11,769 IN records. We observed a total of 3 discordant cases in PR, 3 cases in RT, and none in IN. For cases where both pipelines generated resistance score predictions, we listed the discordant scores for the respective drug names. Putative reasons for discordance were assessed from validation outputs.

Country/region	Sample size	Subtypes	Sequence length (nt)	Accession numbers	Ref.
Brazil	103	B (100%)	1262.0	MF545238 – MF545340	
Ethiopia	67	C (97.0%), B (1.5%), A (1.5%)	1042.1	MH324937 – MH325003	[20]
Guinea-Bissau	54	CRF02_AG (88.9%), A (5.6%), CRF06_cpx (1.8%)	1035.0	MH605452 – MH605505	[21]
Hong Kong	284	C (36.0%), CRF07_BC (36.0%), CRF02_AG (8.8%)	1157.8	MH757122 – MH757405	
South Africa	212	C (100%)	1195.0	MH920641 – MH920852	[22]
Tajikistan	146	A (97.3%), CRF02_AG (2.0%), CRF63_02A1 (0.7%)	1351.1	MH543115 – MH543260	
Uganda	140	D (99.3%), CRF10_CD (0.7%)	864.0	MH925538 – MH925677	

**Table 2** Characteristics of HIV PR-RT population data sets. Sequence data were obtained for an arbitrary selection of recent studies with HIV-1 *pol* sequences deposited in Genbank that spanned both PR and RT. We processed the sequences through both *SierraPy* and *sierra-local* to confirm that the pipelines obtained identical resistance scores. Subtypes listed were obtained from the *SierraPy* pipeline; the subtype classifications obtained from *sierra-local* were highly concordant (94.5% identical). CRF = circulating recombinant form.

285 **Increased Performance over HIVdb Sierra.** Performance and hence, the speed, of software pack-  
286 ages varies according to hardware, software, and input data characteristics. All development, test-  
287 ing, and validation was performed on a workstation running Ubuntu 18.04 LTS with an Intel Xeon  
288 E5-1650 v4 hexa-core CPU at 3.60 GHz and 16 GB of DDR4-2400 RAM with a gigabit net-  
289 work connection. *sierra-local* achieved mean [range] processing speeds of 47.08 sequences/second  
290 (seq/s) [45.07, 48.49] for PR, 16.20 seq/s [14.01, 19.97] for RT, and 14.99 seq/s [14.79, 15.56] for  
291 IN. A substantial fraction of processing time was consumed by subtyping. *SierraPy*, with the same  
292 dataset as previously described, yielded mean processing speeds of 16.01 seq/s [12.88, 17.60] for  
293 PR, 6.12 seq/s [4.83, 7.54] for RT, and 5.19 seq/s [5.05, 5.47] for IN. Although the size of sequence  
294 batches used in this performance comparison likely is a factor in the results by virtue of file writing  
295 and reading being done once per batch, the large batch size used minimizes the effect of these I/O

296 processes on the overall runtime. Overall, *sierra-local* is able to process and return results for sub-  
297 mitted query HIV-1 *pol* sequences roughly 3 times faster than SierraPy, depending on the nature of  
298 the sequences and the type of local computing resources available. This result is in the expected di-  
299 rection since local computing resources are able to be fully utilized, whereas SierraPy depends on  
300 server-side processing speed, server load, request balancing, as well as network speed and traffic.  
301 In this case, network speeds between SierraPy on the local workstation and HIVdb Sierra GraphQL  
302 Web Service were likely not a significant factor in the speed improvement results obtained. With  
303 slower network speeds and all other factors being equal, however, the relative processing speed of  
304 *sierra-local* to SierraPy can only be expected to increase.

305 **Concluding remarks.** The distribution of the HIVdb resistance genotyping algorithm in a stan-  
306 dardized format (ASI [12] is an important resource for HIV-1 research and clinical management,  
307 and an exemplary case of open science. *sierra-local* provides a convenient framework to generate  
308 HIV drug resistance predictions from ASI releases in a secure environment and confers full control  
309 over data provenance. The ability to apply ASI-encoded algorithms locally (offline) also makes  
310 this part of the laboratory workflow robust to network availability may be particularly important  
311 for laboratories situated in resource-limited settings. In addition, the relative processing speed  
312 of *sierra-local* can confer an advantage for research applications requiring the analysis of large  
313 numbers of sequences. The emulation of an established genotype interpretation system to process  
314 unaligned nucleotide sequences and produce identical resistance predictions and data summaries  
315 in a small, standalone package was not a trivial undertaking. Despite the relative simplicity of the  
316 rules-based HIVdb algorithm, there were a large number of pre-processing and post-processing  
317 steps that were necessary to adapt to maximize concordance with the original system. We de-  
318 veloped *sierra-local* with the aims of minimizing the number of additional programs that users  
319 would have to install for a local implementation. The only other presently available means for  
320 implementing for a completely independent instance of the HIVdb algorithm is by hosting an in-  
321 stance of the HIVdb Sierra Web Service itself, which was recently made possible with the release  
322 of the Sierra source code. This approach, however, requires the configuration of a web server,

323 the Apache Tomcat web container, and a large number of Java libraries (Apache Commons Lang,  
324 Apache Commons Math, Apache Commons IO, Apache Log4j, Google Guava, Google Gson, pro-  
325 tonpack, and GraphQL-Java). Furthermore, hosting a stable web service for the sole purpose of  
326 independently generating clinical resistance predictions increases facility requirements for servers,  
327 information technology support staff, and generally complicates an already complex workflow. We  
328 hope that making this lightweight, open-source implementation of the HIVdb ASI to the clinical  
329 and research community will further democratize HIV drug resistance genotyping across providers  
330 of HIV care.

## ACKNOWLEDGEMENTS

331 We thank Philip Tzou for bringing NucAmino to our attention, and for his contributions to open  
332 science in the release of the Stanford HIVdb resistance program source code. This work was sup-  
333 ported in part by the Government of Canada through Genome Canada and the Ontario Genomics  
334 Institute (OGI-131) and by a grant from the Canadian Institutes of Health Research (PJT-156178).  
335 The funders had no role in study design, data collection and interpretation, or the decision to submit  
336 the work for publication.

## REFERENCES

337 [1] **Tural C, Ruiz L, Holtzer C, Schapiro J, Viciana P, González J, Domingo P, Boucher C, Rey-Joly C, Clotet B, et al.** 2002. Clinical utility of HIV-1 genotyping and expert advice:  
338 the Havana trial. *AIDS* **16**:209–218.

340 [2] **Günthard HF, Aberg JA, Eron JJ, Hoy JF, Telenti A, Benson CA, Burger DM, Cahn P, Gallant JE, Glesby MJ, Reiss P, Saag MS, Thomas DL, Jacobsen DM, Volberding PA, International Antiviral Society-USA Panel.** July 2014. Antiretroviral Treatment of Adult  
341 HIV Infection: 2014 Recommendations of the International Antiviral Society-USA Panel.  
342 *JAMA* **312**:410–25. doi:10.1001/jama.2014.8722.

345 [3] **Mayer KH, Hanna GJ, D'Aquila RT.** March 2001. Clinical Use of Genotypic and Phe-

346 notypic Drug Resistance Testing to Monitor Antiretroviral Chemotherapy. *Clin Infect Dis*  
347 **32**:774–782. doi:10.1086/319231.

348 [4] **Liu TF, Shafer RW**. June 2006. Web Resources for HIV Type 1 Genotypic-Resistance Test  
349 Interpretation. *Clinical Infectious Diseases* **42**:1608–1618. doi:10.1086/503914.

350 [5] **Meynard JL, Vray M, Morand-Joubert L, Race E, Descamps D, Peytavin G, Matheron  
351 S, Lamotte C, Guiramand S, Costagliola D, et al.**. 2002. Phenotypic or genotypic resistance  
352 testing for choosing antiretroviral therapy after treatment failure: a randomized trial. *AIDS*  
353 **16**:727–736.

354 [6] **Van Laethem K, De Luca A, Antinori A, Cingolani A, Perno CF, Vandamme AM**. 2002.  
355 A genotypic drug resistance interpretation algorithm that significantly predicts therapy re-  
356 sponse in HIV-1-infected patients. *Antivir Ther* **7**:123–129.

357 [7] **Tang MW, Liu TF, Shafer RW**. 2012. The HIVdb System for HIV-1 Genotypic Resistance  
358 Interpretation. *Intervirology* **55**:98–101. doi:10.1159/000331998.

359 [8] **Bernard EJ, Azad Y, Vandamme AM, Weait M, Geretti AM**. September 2007. HIV  
360 forensics: pitfalls and acceptable standards in the use of phylogenetic analysis as evidence  
361 in criminal investigations of HIV transmission. *HIV Med* **8**:382–387. doi:10.1111/j.1468-  
362 1293.2007.00486.x.

363 [9] **Patil HK, Seshadri R**. 2014. Big data security and privacy issues in healthcare. In *Big Data*  
364 (BigData Congress), 2014 IEEE International Congress on. pages 762–765. IEEE.

365 [10] **Hart SAS, Vardhanabhuti S, Strobino SA, Harrison LLJ**. June 2018. The Impact of  
366 Changes over Time in the Stanford University Genotypic Resistance Interpretation Algo-  
367 rithm. *J Acquir Immune Defic Syndr* **79**:1. doi:10.1097/QAI.0000000000001776.

368 [11] **Aziz N, Zhao Q, Bry L, Driscoll DK, Funke B, Gibson JS, Grody WW, Hegde MR,  
369 Hoeltge GA, Leonard DG, Merker JD, Nagarajan R, Palicki LA, Robetorye RS, Schri-**

370 **Jver I, Week KE, Voelkerding KV.** April 2015. College of American Pathologists' Lab-  
371 oratory Standards for Next-Generation Sequencing Clinical Tests. *Arch Pathol Lab Med*  
372 **139**:481–493. doi:10.5858/arpa.2014-0250-CP.

373 [12] **Betts BJ, Shafer RW.** June 2003. Algorithm Specification Interface for Human Im-  
374 munodeficiency Virus Type 1 Genotypic Interpretation. *J Clin Microbiol* **41**:2792–2794.  
375 doi:10.1128/JCM.41.6.2792-2794.2003.

376 [13] **Shafer RW.** 2006. Rationale and uses of a public hiv drug-resistance database. *The Journal*  
377 of infectious diseases **194**:S51–S58.

378 [14] **Tzou PL, Huang X, Shafer RW.** 2017. Nucamino: a nucleotide to amino acid alignment  
379 optimized for virus gene sequences. *BMC Bioinformatics* **18**:138. doi:10.1186/s12859-017-  
380 1555-6.

381 [15] **Mangeat B, Turelli P, Caron G, Friedli M, Perrin L, Trono D.** July 2003. Broad antiretro-  
382 viral defence by human APOBEC3G through lethal editing of nascent reverse transcripts.  
383 *Nature* **424**:99–103. doi:10.1038/nature01709.

384 [16] **Snoeck J, Kantor R, Shafer RW, Van Laethem K, Deforche K, Carvalho AP, Wynhoven**  
385 **B, Soares MA, Cane P, Clarke J, Pillay C, Sirivichayakul S, Ariyoshi K, Holguin A,**  
386 **Rudich H, Rodrigues R, Bouzas MB, Brun-Vezinet F, Reid C, Cahn P, Brigido LF,**  
387 **Grossman Z, Soriano V, Sugiura W, Phanuphak P, Morris L, Weber J, Pillay D, Tanuri**  
388 **A, Harrigan RP, Camacho R, Schapiro JM, Katzenstein D, Vandamme AM.** February  
389 2006. Discordances between Interpretation Algorithms for Genotypic Resistance to Protease  
390 and Reverse Transcriptase Inhibitors of Human Immunodeficiency Virus Are Subtype De-  
391 pendent. *Antimicrob Agents Chemother* **50**:694–701. doi:10.1128/AAC.50.2.694-701.2006.

392 [17] **Wagner S, Kurz M, Klimkait T.** 2015. Algorithm evolution for drug resistance  
393 prediction: comparison of systems for HIV-1 genotyping. *Antivir Ther* **20**:661–665.  
394 doi:10.3851/IMP2947.

395 [18] **R Core Team.** 2018. R: A Language and Environment for Statistical Computing. R Founda-  
396 tion for Statistical Computing. Vienna, Austria.

397 [19] **Ooms J.** 2014. The jsonlite Package: A Practical and Consistent Mapping Between JSON  
398 Data and R Objects. arXiv:1403.2805 [stat.CO] .

399 [20] **Arimide DA, Abebe A, Kebede Y, Adugna F, Tilahun T, Kassa D, Assefa Y, Balcha TT, Björkman P, Medstrand P.** 2018. Hiv-genetic diversity and drug resistance transmission  
400 clusters in gondar, northern ethiopia, 2003-2013. PloS one **13**:e0205446.

402 [21] **Wilhelmsen S, Måansson F, Lindman JL, Biai A, Esbjörnsson J, Norrgren H, Jansson M, Medstrand P, et al.** 2018. Prevalence of hiv-1 pretreatment drug resistance among treatment  
403 naïve pregnant women in bissau, guinea bissau. PloS one **13**:e0206406.

405 [22] **Rasmussen D, Wilkinson E, Vandormael A, Tanser F, Pillay D, Stadler T, de Oliveira  
406 T.** 2017. External introductions helped drive and sustain the high incidence of hiv-1 in rural  
407 kwazulu-natal, south africa. bioRxiv page 119826.