

1 **STRyper: a macOS application for microsatellite genotyping and chromatogram**
2 **management**

3 Jean Peccoud^{1*}

4

5 ¹Laboratoire Écologie et Biologie des Interactions, Équipe Écologie Évolution Symbiose,
6 Université de Poitiers, UMR CNRS 7267, Poitiers, France

7 * Corresponding author

8 Email: jeanpeccoud@gmail.com

9

10

11

12 Abstract

13 Microsatellite markers analyzed by capillary sequencing remain useful tools for rapid
14 genotyping and low-cost studies. This contrasts with the lack of a free application to analyze
15 chromatograms for microsatellite genotyping that is not restricted to human genotyping. To
16 fill this gap, I have developed STRyper, a macOS application whose source code is published
17 under the General Public License. STRyper only uses macOS libraries, making it very
18 lightweight, responsive, and behaving like a modern application. Its three-pane window
19 enables easy management and viewing of chromatograms imported from .fsa and .hid files,
20 the creation of size standards and of microsatellite marker panels (including bins). STRyper
21 features powerful search capabilities (with smart folders) and a modern graphical user
22 interface allowing, among others, the manual correction of DNA ladders and of individual
23 genotypes by drag-and-drop. It also introduces a new way to mitigate the effect of variations
24 in electrophoretic conditions on estimated allele sizes.

25

26 Keywords: microsatellites, capillary electrophoresis, chromatograms, population genetics,
27 graphical user interface

28 Introduction

29 More than three decades after their first use, microsatellites markers, also known as short
30 tandem repeat (STR) loci, remain popular DNA markers to assess gene flow, population
31 history, structure and membership, ancestry, or the integrity of laboratory breeding lines,
32 among other uses [1, 2]. When locus-specific variation is not the focus of a study, a limited
33 number of microsatellite markers are sufficient to assess evolutionary processes affecting the
34 whole genome and to genetically identify an individual [3]. This ability stems from the sheer
35 number of alleles per marker, which often counts in the dozens, leading to a per-locus
36 information amount that exceeds that of single-nucleotide polymorphisms (SNPs) [4].

37 Due to frequent indels affecting the number of microsatellite repeat motives, microsatellites
38 alleles essentially differ in their length, which can be estimated by simple electrophoresis of
39 amplicons. Amplicon sequencing by the Illumina technology has however emerged as

40 relatively affordable and more reliable alternative to capillary electrophoresis De Barba,
41 Miquel (5), Barbian, Connell (6), Suez, Behdenna (7)]. In species for which tried and tested
42 microsatellite multiplexes exist, microsatellite genotyping via electrophoresis still offers a
43 compelling money- and time-saving solution. At a few dollars per individual in terms of
44 consumables (for a couple of multiplexes typically combining 10-20 loci) genotyping can be
45 performed locally in one day, as it amounts to DNA extraction, PCR, amplicon dilution and
46 placing a plate in a capillary sequencer. When a quick answer is needed or when only few
47 individuals need analyzing, typically for simple genotype checking, this traditional technique
48 remains the cheapest and easiest one.

49 However, the difficulty sharply rises when it comes to analyzing the results of capillary
50 electrophoresis. As opposed to genotyping via NGS, which is generally done via fully- or
51 partially automated free tools (e.g., [8, 9]), traditional microsatellite genotyping requires
52 inspecting fluorescence curves, therefore applications with a complex graphical user interface
53 (GUI), which are rarely free. To various degrees, these applications are focused on human
54 identification by genotyping and forensics. As such, they are packed with features and
55 safeguards that are of little relevance to most researchers, which somewhat complicate their
56 use, and which may come at a high price.

57 This is the case of GeneMapper by ThermoFisher Scientific, a commercial application
58 running on the Windows operating system, and which remains, to my knowledge, the most
59 widely used for microsatellite genotyping. A Google scholar search for “genemapper”,
60 excluding references, patents and review articles, and limited to 2023 and 2024, returned 2820
61 results as of July 20th 2024. Most results pertained to medicine and forensics or may
62 correspond to preprints, but the first 130 results comprised 15 English-written studies on non-
63 human species using traditional microsatellite analyzes, indicating that this technique is far
64 from abandoned.

65 GeneMarker by Softgenetics is a similar commercial application. The price of a license of
66 either software may restrict its installation to a single computer per research laboratory. A free
67 alternative from ThermoFisher Scientific, Peak Scanner, has limited functionalities.

68 Complementary command-line tools [10, 11] provide missing features such as allele scoring
69 via binning, but may dissuade those who seek to conduct fragment analyses from
70 chromatogram import to the export of individual genotypes in a single user-friendly
71 application. In that regard, Geneious Prime and its microsatellite analysis plugin may
72 represent an interesting tradeoff between price and features. The cost of a subscription to

73 Geneious Prime may still appear excessive to users who do not need the features that this
74 product offers for the analysis of DNA sequences.

75 Osiris [12, 13], stands out as being a free, feature-rich and multi-platform (Windows and
76 macOS) tool for STR analysis. Yet, this software is, as far as I know, rarely used by
77 population geneticists, possibly because it is highly specialized for human identification.

78 Researchers, especially population geneticists, would therefore benefit from a free application
79 enabling quick microsatellite genotyping and management of thousands of samples. To meet
80 this need, I have developed STRyper, an open-source, lightweight and user-friendly
81 application that can analyze chromatogram files for STR genotyping. STRyper is published
82 under the GNU General Public License v. 3 and its name is a portmanteau of “STR” and
83 “Genotyper”. As described below, STRyper features a modern GUI allowing, among others,
84 unconstrained chromatogram management via nested folders, advanced and dynamic
85 metadata-based chromatogram search with “smart” folders, easy folder import/export,
86 chromatogram and genotype filtering based on multiple criteria, the definition of
87 microsatellite multiplexes and custom size standards, fast and responsive visualization of
88 fluorescence curves with animated zooming and automatic vertical scaling, the manual
89 correction of DNA ladders and of individual genotypes by drag-and-drop, and a new way to
90 mitigate the effect of variations in electrophoretic conditions on estimated allele sizes. The
91 application and its codebase are available at <https://github.com/jeanlain/STRyper>.

92 Description of the application

93 General characteristics and development

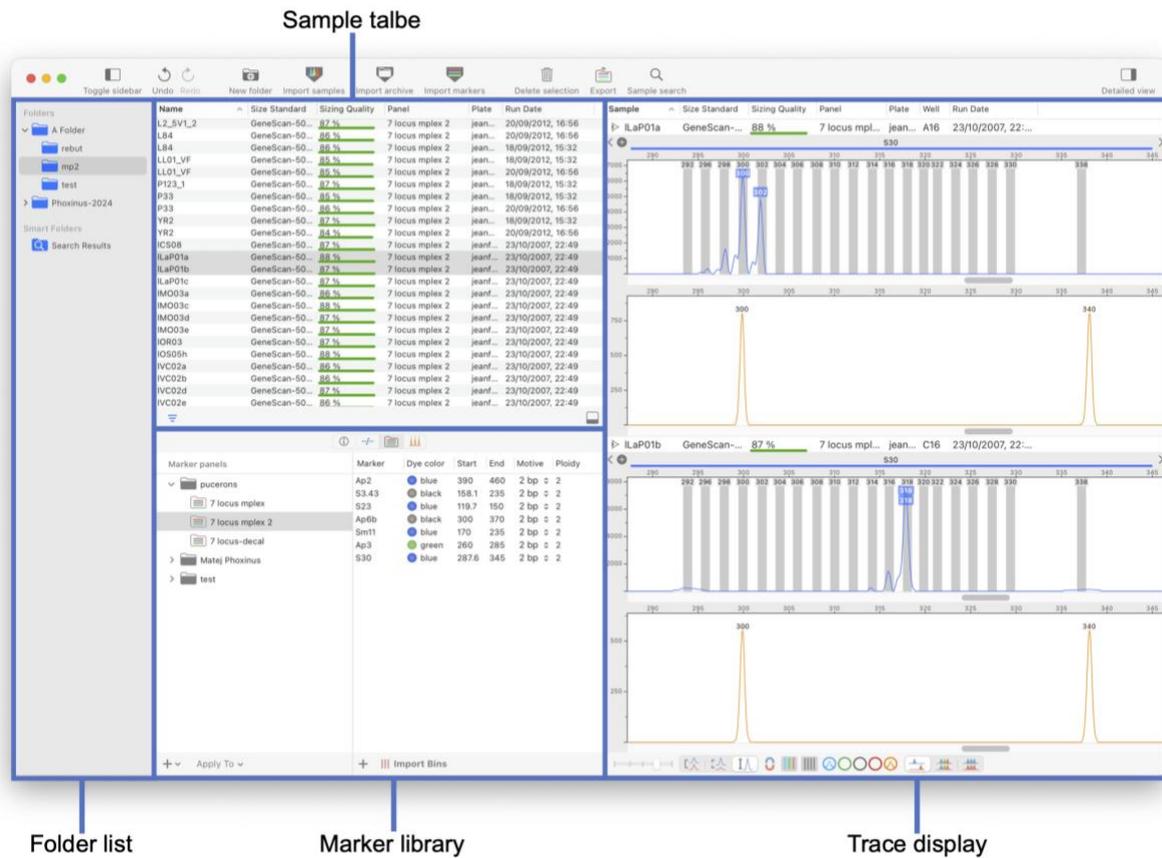
94 STRyper is designed to manage and display chromatograms generated by Applied Biosystems
95 capillary sequencers. The application also allows managing microsatellite markers and size
96 standards, which are required for fragment size estimation and genotyping. All functionalities
97 of the application are driven by its GUI. As opposed to command line tools, GUI development
98 relies on application programming interfaces and frameworks that depend on the target
99 operating system and development tools. These were dictated by my use of the Mac operating
100 system (macOS) and by the fact that developing STRyper was a hobby project of an
101 evolutionary biologist, not the effort of a team of professional developers. Being
102 unencumbered by cross-platform development gave me the freedom to choose the right tools

103 to program a GUI that was intuitive, responsive and consistent with “native” macOS
104 applications. STRyper was thus developed using Xcode and frameworks provided by Apple
105 (https://developer.apple.com/library/archive/documentation/MacOSX/Conceptual/OSX_Tech
106 [nology_Overview/SystemFrameworks/SystemFrameworks.html](https://developer.apple.com/library/archive/documentation/MacOSX/Conceptual/OSX_Tech/nology_Overview/SystemFrameworks/SystemFrameworks.html)). These frameworks include
107 “Core Data”, which is used to define and manage objects representing chromatograms,
108 microsatellite marker, bins, alleles, genotypes and size standards, and to save them in a
109 persistent relational database (S1 Text). Internally, Core Data relies on the SQLite database
110 engine to manage the persistent store. GUI elements (windows, views, controls and so on) are
111 implemented using “AppKit”. “Core Graphics” functions are used to draw fluorescent curves.
112 “Core Animation” layers accelerate compositing via the graphical processing unit (GPU) and
113 provide fluid animation of the interface. These object-oriented frameworks (except Core
114 Graphics) required the use of the Objective-C programming language (a superset of C) when
115 the project started. The application code was written in the latest version (2.0) of this
116 language.

117 STRyper runs under macOS version 10.13 or higher. The application does not include third-
118 party libraries and does not require special installation steps. Its bundle contains binaries
119 compiled for the X86 and arm64 architectures and weighs less than 15 Megabytes, including
120 the user guide.

121 Overview of the interface

122 The application comprises a main window (**Fig 1**) composed of three panes, a design
123 paradigm used by several database-management applications like email clients. The left
124 collapsible sidebar is a hierarchical list of folders and subfolders containing samples, each
125 representing an imported chromatogram file. Folder and samples can be organized freely by
126 drag and drop. A middle pane shows the content of the selected folder (samples and
127 associated genotypes) and comprises tabs to manage size standards and markers. The right
128 pane shows the traces (fluorescent curves) of selected samples and genotypes.



130 **Fig 1. The main window of STRyper.** The left pane contains the list of folders and smart
 131 folders (search results) containing samples. The middle pane is a split view comprising a top
 132 pane listing the samples of the selected folder. Its bottom pane has four tabs, which are from
 133 left to right: an inspector showing data on selected samples (Fig 2), a table of genotypes from
 134 the samples shown on the top pane, the marker library (currently shown) and the size standard
 135 library. The right pane shows the traces of selected samples in a scrollable view that can
 136 display thousands of traces.

137 STRyper uses very few modal panels or dialogs to validate user actions and all actions that
 138 affect the database can be undone. Most can be achieved in a couple of clicks or less as they
 139 do not require opening and closing windows. Drag and drop can be used throughout: from
 140 importing samples to applying size standards, markers, and to manually attributing alleles or
 141 size molecular ladder fragments to peaks.

142 STRyper can import FSA files (HID file support is experimental, as the HID format
 143 specifications are not public) containing data for 4 or 5 channels (fluorescent dyes). Samples
 144 are imported into folders, and they can be moved or copied between folders at any time. A
 145 folder and all its content, including subfolders, samples, genotypes at microsatellite markers,

146 associated marker panels (including bins) and custom size standards, can be archived and
147 transferred between instances of the application. Upon importing an archived folder, any
148 marker panel (multiplexes) and size standard encoded in the archive is imported unless it
149 already in the database. The imported folder therefore shows the same content as the original
150 one.

151 Since samples are not constrained to compartmentalized projects, the application provides
152 search tools to find and gather samples from the whole database. Users can define various
153 search criteria, including run date, sizing quality, well identifier, plate name, marker panel
154 name, etc. Search results appear in “smart folders” which dynamically update their contents as
155 new samples meet the search criteria.

156 Chromatogram display

157 Like all applications displaying chromatograms generated by capillary sequencers, STRyper
158 draws plots in which the Y-axis is the fluorescence level. The X-axis represents the length of
159 DNA fragments that produced peaks in the fluorescence. This contrasts with Osiris, in which
160 the X-axis is the number of the fluorescence data record (the “scan” number), hence the time
161 at which the measure was taken during the electrophoresis.

162 For fragment size estimates, the application must first identify fluorescence peaks, which are
163 induced by DNA fragments. This task is performed during chromatogram import by a simple
164 algorithm. This algorithm (detailed in the S1 Text) determines whether the fluorescence level
165 at a given scan is elevated enough, both relative to neighbor scans and in absolute level. Peak
166 delineation serves as a basis to subtract baseline fluorescence level, which helps peak
167 visualization. The method developed for this task adjusts the height (fluorescence level) of a
168 curve such that the start and end point of each peak are placed at level zero (S1 text).

169 Although this adjustment cannot be applied on signals that are too faint to contain meaningful
170 peaks, it has the benefit of offering two baseline subtraction modes: one that preserves
171 absolute peak height, and one that maintains relative peak elevation compared to the baseline
172 (S1 Text). As this method reduces background noise, no smoothing algorithm was
173 implemented.

174 Because chromatograms contain fluorescence data from several wavelengths (channels),
175 multichannel fluorescence analysis requires determining whether a peak represents a DNA
176 fragment or interference from another channel (i.e., “crosstalk”). The method developed for
177 this task compares the position, shape and relative size of peaks between channels, accounting

178 for saturation of the sequencer camera (S1 Text). To signal crosstalk to the user, the area
179 underneath an artefactual peak is filled with the color that represents the channel that induced
180 crosstalk (this option can be disabled). While certain applications alter fluorescence data to
181 correct for pull-up due to crosstalk [13], flagging peaks resulting from crosstalk and leaving
182 the source signal untouched was considered sufficient. These peaks are simply ignored in
183 automatic detection of alleles and DNA ladder fragments (detailed below), although the user
184 can manually assign these peaks, should they wish to.

185 Upon selecting samples in the table, corresponding fluorescent curves (traces) are
186 instantaneously displayed on the right pane (**Fig 1**). As the application fully supports the dark
187 theme of macOS (version 10.14 or more recent), it can display traces on a dark background to
188 alleviate eye strain. Any region in which a peak saturated the sequencer camera is shown
189 behind curves as a rectangle whose color reflects the channel that likely caused saturation.
190 Traces can be scrolled and zoomed in/out horizontally via trackpad gestures such as swipe,
191 pinch and double tap, via the scroll wheel, or by clicking/dragging the mouse over horizontal
192 rulers to define a size range. Dragging the mouse over the vertical ruler sets the fluorescence
193 level at the top of the view, hence the vertical scale. Zooming is animated, which helps users
194 keep track of the range (in base pairs) that is displayed.

195 Viewing options include automatic vertical scaling to the highest visible peaks, synchronizing
196 the vertical scales and horizontal positions, showing/hiding region of fluorescence saturation,
197 stacking curves from several samples or channels in the same view, and subtracting the
198 baseline fluorescence level.

199 Size standards and molecular ladders

200 To estimate the size of DNA fragments, a molecular ladder containing fragments of know
201 lengths (defining a “size standard”), and tagged with a specific fluorescent dye, is added to
202 every sample before electrophoresis. DNA ladder fragments induce peaks in the trace of the
203 corresponding channel. To associate peak to sizes, samples must be assigned the adequate size
204 standard. STRyper comes with several widely used size standards, namely those from the
205 GeneScan brand. Users can easily edit these size standards within the application and make
206 their own. They can be assigned to samples during chromatogram import (based on metadata
207 encoded in the file) or manually. Assigning a size standard automatically triggers the
208 detection of DNA ladder fragments in the sample.

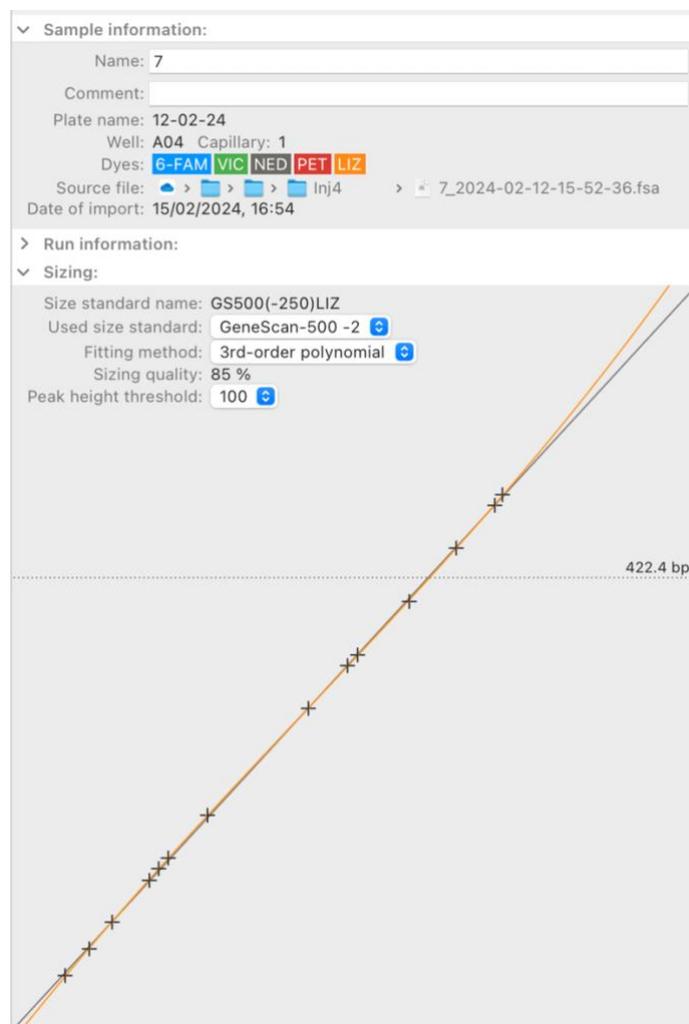
209 The method used to detect DNA ladder fragments and assign them to sizes of a known size
210 standard is based on relative peak positions and accounts for non-linear relationship between
211 fragment size and migration speed (S1 Text). Peaks resulting from crosstalk or whose height
212 are unusual compared to others are ignored. To account for non-linearity, a polynomial of the
213 first, second, or third degree (depending on the user choice) is used to estimate fragment size,
214 where the response variable is the size of a fragment specified in the size standard, and the
215 explanatory variable is the scan numbers at the tip of the corresponding peak (representing
216 migration speed). This principle is also implemented in other applications such as
217 GeneMapper. Fitting is achieved via the Cholesky decomposition implemented in the Linear
218 Algebra Package (<https://netlib.org/lapack/>). Fitting parameters are used to draw traces by
219 computing the size in base pairs corresponding to every scan. The horizontal distance between
220 successive scans varies unless a polynomial of the first degree (linear regression) is used for
221 the sizing.

222 To evaluate the quality of the sizing, a score from 0 to 1 was developed, based on the
223 residuals of the fitted model (differences between fragment sizes as defined in the size
224 standard, and fragment sizes estimated by the model). This score involves computing the
225 difference in residuals for every pair of adjacent peaks and is computed as follows. If ΔR is
226 the difference between residuals of every pair of adjacent peaks, ΔS the difference in scan
227 number of these peaks, n_p the number of peaks and n_s is number of sizes in the size standard,
228 the quality score is:

$$229 \quad 1 - \max\left(\frac{\Delta R^2}{|\Delta S|}\right) \frac{10}{3} - \frac{n_s - n_p}{10}$$

230 Any negative score is set to zero. This formula was tuned by testing many chromatograms.
231 The 10/3 coefficient ensures that the score is greatly reduced (often to zero) by a single
232 assignment error, which affects $\max\left(\frac{\Delta R^2}{|\Delta S|}\right)$. A poor score compels the user to rectify the error.
233 The score is also reduced if certain sizes of the size standard are not assigned to any peak,
234 which increases $n_s - n_p$. A weight of 1/10 was attributed to this component, because such issue
235 generally reflects problems during electrophoresis, which cannot be fixed in the application.
236 Sizing quality is shown for each sample in a dedicated column displaying a gauge (**Fig 1**). If
237 molecular sizing failed, sizes are not displayed on the X-axis of the chromatogram, but traces
238 can still be viewed.

239 STRyper displays the trace of the molecular ladder like any other trace, letting users switch
240 quickly between genotype and molecular ladder editing. Sizes attributed to molecular ladder
241 fragment can be changed by dragging and dropping size labels onto peaks. Any change to the
242 molecular ladder automatically updates the sizing of the sample without user validation. The
243 red component of the color used for size labels is proportional to the difference between the
244 computed size of a peak and its theoretical size, making size assignment errors easy to spot.
245 The application also features an inspector panel that dynamically updates to show metadata of
246 selected samples, and most importantly, sizing information (**Fig 2**~~Error!~~ **Reference source**
247 **not found.**). This inspector can help to find sizing errors if points deviate from the curve
248 representing the relationship between scan number and peak size.



249
250 **Fig 2. The sample inspector of STRyper.** This panel with three collapsible sections
251 dynamically updates to display information on samples that are selected in the sample table
252 (**Fig 1**). The plot at the bottom shows the relationship between the time at which DNA
253 fragments of the molecular ladder (black crosses) were detected by the sequencer camera (the

254 X axis) and their observed sizes in base pairs (the Y axis). The relationship used to estimate
255 fragment sizes is established by fitting a polynomial (here, of the third degree) to the points
256 shown on the plot. This polynomial is represented by the orange curve. The horizontal dotted
257 line indicates the estimated size at the location of the mouse cursor (cursor not shown).

258 **Microsatellite marker and bins**

259 Genotyping requires associating chromatograms with the microsatellite markers that were
260 amplified using fluorescent primers. Markers amplified together by multiplex PCR are
261 regrouped into a “panel” (a term derived from GeneMapper). Users can define their own
262 panels of haploid or diploid microsatellite markers within STRyper and organize them in
263 folders. Markers are defined by their fluorescent dye, ploidy, length of repeat motive, name,
264 and the size range of their alleles. These attributes can be changed after a marker is created,
265 except for the first two. Markers can be copied and dragged between panels. Users can export
266 marker panels to text files conforming to simple specifications described in the user guide.
267 These text files can be imported back as marker panels. STRyper can also import panel
268 description text files exported from GeneMapper.

269 A marker can comprise “bins”, which are non-contiguous intervals delimiting the expected
270 sizes of fragments corresponding to alleles [14]. Bins address the fact that estimated fragment
271 sizes slightly vary between sequencer runs [15]. Proper bin definition must account for factors
272 affecting amplicon mobility during electrophoresis [16], which often cause the estimated
273 distance between consecutive microsatellite alleles (in base pairs) to slightly differ from the
274 repeat motive length [14]. Binning can be left to specialized programs like Tandem [17],
275 which can work on allele sizes estimated by other programs like STRyper. The management
276 of bins within STRyper was still considered a necessity. Indeed, visualizing bins as vertical
277 rectangles behind traces helps to characterize alleles that do not conform to the periodicity of
278 the repeat motive, and to mitigate variations in fragment sizes between sequencer runs
279 (further discussed below). STRyper therefore allows importing bin sets as text files (produced
280 by GeneMapper or Tandem), but also generating and editing bin sets within the application.

281 In STRyper, a set of automatically named bins for a marker can be added by specifying the
282 width and spacing of bins. To accommodate the fact that the observed distance between
283 microsatellite alleles slightly differs from the repeat motive [14], the position and a spacing of
284 bins can be adjusted by respectively dragging and resizing the whole bin set. Individual bins
285 can also be added and modified via click and drag. These actions do not involve dedicated

286 windows or panels, they can be performed at any time on the trace views where bins are
287 displayed (**Fig 1**, right).

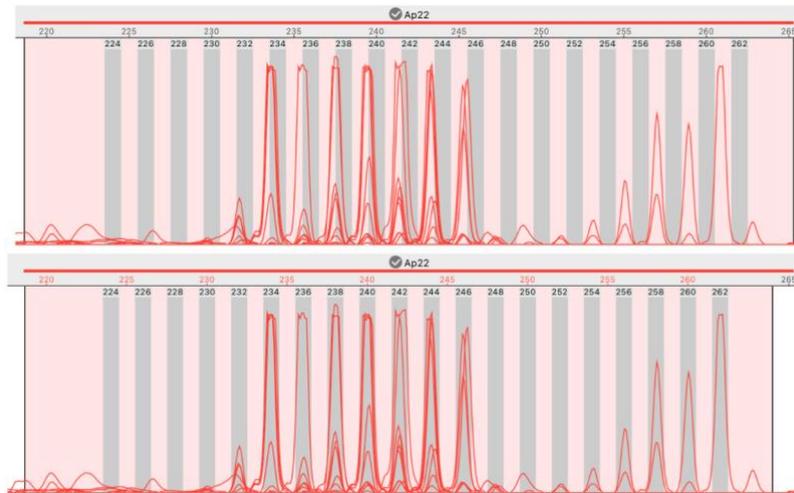
288 The width of a bin might not cover the full range of estimated sizes of amplicons from given
289 allele over all electrophoretic conditions. Regularly, a peak representing an allele would fall
290 outside the corresponding bin, although identical fragments that migrated in other sequencer
291 runs were properly binned. To circumvent the issue, a mixture of amplicons of known sizes
292 for each marker, known as “allelic ladder” or “inter-lane standard”, can be added alongside
293 samples for each run or sequencing plate. Allelic ladders are however only available for
294 model species.

295 STRyper implements a novel approach to mitigate this issue. Rather than moving bins to
296 match peak positions (which requires maintaining several sets of bins per marker), this
297 approach considers that it is the estimated sizes of peaks (in base pairs), not the position of
298 bins, which should be adjusted. The method thus correct fragment sizes using the formula $y =$
299 $a + bx$, where x is the size of a DNA fragment that is estimated by the DNA ladder via the
300 fitted model mentioned earlier, y is the adjusted size, and a and b are constants (hereafter
301 called “offset parameters”). This approach assumes that the effect of varying electrophoresis
302 conditions can be compensated by this linear combination. If there is no correction, $a = 0$ and
303 $b = 1$. Good offset parameters are those that minimize the distance (in base pairs) between
304 peaks and their corresponding bins. Because automatically determining which bins and peaks
305 to associate might have been error-prone, a manual GUI-based method was developed. The
306 application lets the user move and/or resize a rectangle representing the range of the bin set
307 such that bins coincide with peaks (**Fig 3**). To infer offset parameters a and b from this
308 operation, we let s represents the start of a bin and e its end, in base pairs. If s' and e' represent
309 the corresponding boundaries after the user has moved the bin set appropriately, the offset
310 parameters can be computed by solving

$$311 \quad \begin{cases} s' = a + bs \\ e' = a + be \end{cases}$$

$$312 \quad \text{Hence } b = \frac{e' - s'}{e - s} \text{ and } a = s' - s \frac{e' - s'}{e - s}.$$

313 Since the user moves the bin set as a whole, the operation yields same offset parameters for all
314 bins. These parameters are then associated to the chromatograms involved in the procedure
315 (e.g., those displayed in **Fig 3**) and a given marker.



316

317 **Fig 3. A case of out-of-bin alleles that is solved.** Both images show the stacked traces from 8
 318 samples of the same sequencer run. Peaks represent amplicons of a dinucleotide marker called
 319 “Ap22”. Its range is represented by a horizontal red segment above the ruler showing
 320 graduations in base pairs (bp). Bins appear as grey rectangles. Top: peaks are shifted to the
 321 left with respect to bins, and more so for longer alleles, although bins are separated by exactly
 322 two base pairs. Bottom: the user has moved and narrowed the light-pink rectangle
 323 representing the range of the marker, such that bins coincide with peaks. This move translates
 324 into offset parameters $a = -6.40$ and $b = 1.029$ (see main text). As a result, the estimated size
 325 of peaks overlapping bin 258 (bottom image) has changed from ~ 257 bp to ~ 258.1 bp.

326 Genotyping

327 In STRyper, a panel of microsatellite markers is associated to samples via dragging and
 328 dropping a panel icon onto the sample table (**Fig 1**) or via contextual menus, which initializes
 329 a genotype for each sample and each marker of the panel. Genotypes contain no allele
 330 information until alleles are called. Allele calling can be done manually by clicking peaks
 331 within a marker’s range or automatically, via the implementation of a new algorithm.

332 This algorithm accounts for two main biochemical processes producing DNA fragments of
 333 different lengths. One is the addition of a non-template nucleotide to the 3’ end of the new
 334 DNA strand by the DNA polymerase during PCR [18]. Because the added nucleotide is
 335 generally an adenosine, this process is referred to as “adenylation”. If adenylation affects only
 336 a portion of the amplicons, they may differ in length by one nucleotide, generating two peaks.
 337 The other process is “slippage” during replication, causing indels in the repeated region [19].
 338 Slippage may result in a range of different amplicons that differ by the size of the repeat, a

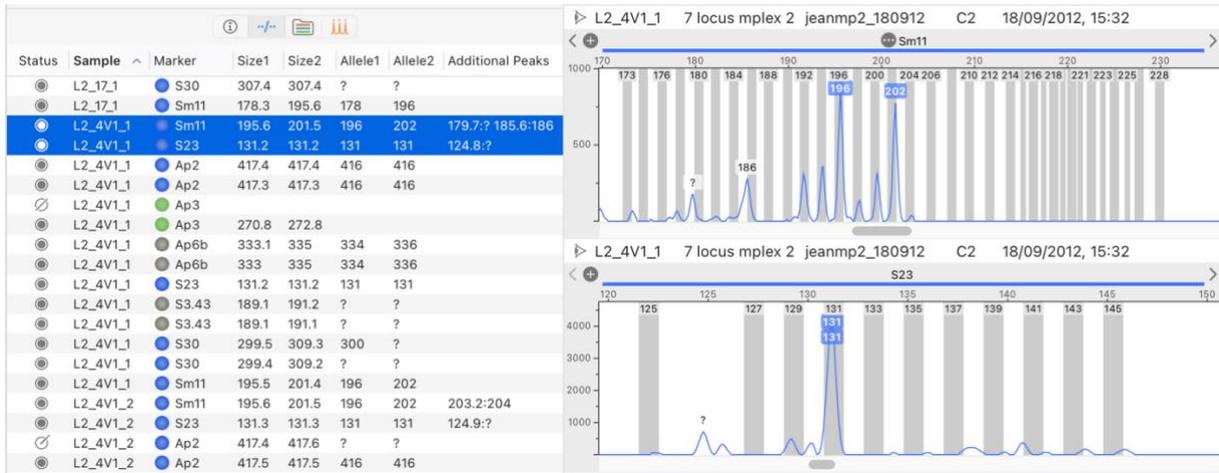
339 pattern known as “stuttering”. These considerations served as a basis to develop a method for
340 allele calling that first identifies peak clusters resulting from these processes (detailed in the
341 S1 Text), and which accounts for the length of the repeat motive. In each delineated cluster,
342 the most intense peak is considered as that representing the allele. Estimation of peak intensity
343 accounts for clipping due to saturation of the fluorescence signal, in that the width of the
344 saturated region is used when peak height/area may not reflect the quantity of DNA material.
345 Stuttering and adenylation are managed internally by the application, the user remains free to
346 manually assign an allele to any peak.

347 Importantly, the method does not consider the absolute height or shape of a peak to call the
348 first allele, beyond the fact that a minimal fluorescence level is required to delineate a peak
349 (see S1 Text). If a peak is detected in the marker range (see below) and is not interpreted as
350 crosstalk, at least one allele will be called. It was considered that the assessment of peak
351 quality was better left to the user, who is expected to visually inspect every genotype.

352 For a diploid individual, the number of different alleles detected within a marker’s range
353 determines the individual’s genotype: homozygous if one allele is detected, heterozygous
354 otherwise. Because this inference is invalid for polyploid markers, it was decided that only
355 haploid and diploid markers could be defined in the application, constraining the maximum
356 number of alleles per locus to 2. To cope with this constraint, the ability to annotate additional
357 DNA fragments of interest, either automatically or manually, was implemented. Additional
358 fragments may inform on the presence of paralogs, polyploidy, insufficient specificity of the
359 PCR or contamination between samples. The application therefore distinguishes two types of
360 peaks: those that are interpreted as alleles and whose number is limited to the ploidy of the
361 marker, and others representing these additional DNA fragments. Because neither should
362 comprise fragments produced by stuttering or adenylation, additional peaks are detected like
363 alleles are (i.e., by identifying peak clusters). The relative height of peaks is used to categorize
364 alleles (higher peaks) and additional peaks (smaller peaks).

365 All genotypes from samples of the current folder are listed in a table (**Fig 4**) that can be sorted
366 and filtered according to various criteria (including allele names and sizes). This table lets
367 users quickly scan genotypes, as corresponding peaks and allele labels of the selected
368 genotype(s) appear on the right-pane. Correcting errors in allele call typically takes a single
369 step: the user can simply drag the mouse from a peak to a bin, drag an allele label from one
370 peak to another (**Fig 5**), or double-click a peak, which removes/attaches an allele from/to the
371 peak. Double clicking allele labels lets users enter arbitrary allele names directly above peaks.

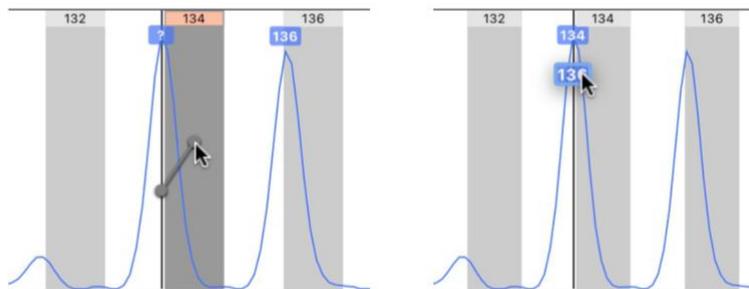
372



373

374 **Fig 4. The table listing genotypes in STRyper.** Each row represents a genotype at a unique
 375 sample and marker. Displayed traces in the right pane correspond to the selected genotypes.
 376 In this instance, two sections of the same trace are shown, corresponding to the range of each
 377 molecular marker. Peaks representing alleles are tagged with labels colored after the channel
 378 of the molecular marker. Labels with a white background represent additional peaks that were
 379 automatically flagged by the allele caller.

380



381

382 **Fig 5. Genotype editing by drag and drop in STRyper.** Vertical grey rectangles represent
 383 bins that define expected ranges of microsatellite alleles. Each bin has a name displayed on
 384 top. Allele names are represented by colored labels above peaks. Left-hand screen capture: the
 385 user is dragging the mouse from a peak to a bin. This will assign the peak an allele named
 386 after the bin, thereby replacing the question mark used for alleles that are out of bins. During
 387 the operation, a grey-colored handle connects the mouse location to another point horizontally
 388 located at the peak tip and vertically located at the clicked point. Right-hand screen capture:
 389 the user has decided that only the peak on the left should represent an allele and is dragging
 390 an allele label from the right-hand peak to the other. These actions are assisted by

391 “magnetism” to lock the handle or allele label to the closest suitable destination, which
392 triggers haptic feedback on the trackpad.

393 Exporting results

394 STRyper allows exporting results in several ways. Genotypes and associated sample metadata
395 can be exported as text files, or simply copied from selected table rows to a text editor or a
396 spreadsheet application. In addition, a folder or a smart folder, with all its content –
397 subfolders, samples, genotypes at microsatellite markers, associated marker panels (including
398 bins) and custom size standards – can be archived and transferred between instances of the
399 application. Upon importing an archived folder, any marker panel and size standard encoded
400 in the archive is imported unless is it already in the database. The imported folder therefore
401 shows the same content as the original.

402 Evaluation of the software

403 Usage and reliability

404 STRyper was developed to facilitate the genotyping of numerous individuals from
405 chromatogram import, management and viewing, to genotype editing and data export. How
406 well it performs at these tasks cannot be evaluated without subjectivity.

407 The ability of an application to assign the right peaks to a DNA ladder fragments or alleles
408 (i.e., allele calling) can be quantified more objectively by comparing these assignments to a
409 reference, which is the assignments that an experienced user would have made by visually
410 inspecting the chromatograms. Another reference, which could be used to evaluate the allele
411 caller specifically, is the genotypes obtained at the same markers from an independent and
412 more reliable method, typically amplicon sequencing. Such reference would allow detecting
413 errors that even an experienced user would not detect. These errors may arise from variations
414 in the motility of amplicons (leading to migration speed not being proportional to fragment
415 length), due to the intrinsic properties of these fragments or variations in experimental
416 conditions, including instruments and operators (reviewed in [20]). Mitigating errors that are
417 visually undetectable should not reasonably be expected from this application. Comparing
418 genotypes called by STRyper to those obtained by sequencing would therefore not
419 constitute a fair evaluation of the allele caller, even if sequence data were available for the
420 same individuals and markers (I am not aware of the public availability of such dataset). The

421 frequency of manual corrections that an experienced user must apply to automatic peak
422 assignment was therefore used as a metric of the application performance, even though it is
423 partly user dependent.

424 Given these limitations, it was considered more valuable to evaluate STRyper as part of an
425 ongoing study (Vucić et al., in prep) instead of reanalyzing previously published data.
426 Chromatograms were obtained from 314 individuals of the freshwater fish *Phoxinus*
427 *lumaireul* (Teleostei, Cypriniformes), each amplified at two 6-plexes of microsatellite
428 markers developed by Vucic, Jelic (21). Amplicons were submitted to electrophoresis in an
429 SeqStudio sequencer (Applied Biosystems) after addition of the GeneScan 500-LIZ size
430 standard. After importing the 648 chromatograms (628 from amplified samples and 20
431 negative controls) into STRyper, the GeneScan 500 size standard was applied to each using
432 the 3rd degree polynomial as sizing method.

433 Ignoring electrophoresis failures that made 25 samples unusable, visual inspection of peak
434 assignments to DNA ladder fragments revealed issues in eight chromatograms. In all cases, a
435 size was not assigned to the appropriate peak or the DNA ladder because the peak was
436 missing or abnormally short. As issues due to missing peaks cannot be fixed, manual
437 corrections were applied to only four chromatograms. Overall, the verification of the DNA
438 ladder for all chromatograms took less than five minutes.

439 For each marker, a set of bins was generated in one step by a specifying a bin width of 1 base
440 pair and setting bin spacing according to the length of the microsatellite repeat motives [21].
441 For samples of a reference sequencing plate, the bin set was moved and resized as a whole,
442 such that bins position matched peaks corresponding to alleles. For certain other sequencing
443 plates and for five markers (PHOX4, PHOX11, PHOX29, PHOX33 and CtoA-247 [21]),
444 peaks and bins appeared slightly misaligned (by less than 0.5 base pairs). Offset parameters
445 for peak sizes were thus defined according to the procedure shown in **Fig 3**. This procedure
446 made bins coincide neatly with peak locations for all regularly spaced alleles. I therefore saw
447 no evidence that the use of linear relationship to estimate offset parameters was inappropriate.
448 Individual bins were also added at locations indicating the presence of alleles that did not
449 strictly follow the repeat pattern (probably due to mutations in microsatellite flanking
450 regions).

451 Once these adjustments were done, genotypes were called and visually checked. Marker
452 PHOX02 suffered from a combination of high stuttering, variable adenylation rates, the
453 probable existence of mutations in flanking regions, which made peak assignment and binning

454 very difficult, even visually. The marker was excluded, because it was considered too
455 unreliable.

456 In several cases of PCR failures, the application assigned relatively faint peaks amounting to
457 noise as alleles, which was expected. These cases were easily detected by visual inspection.
458 The only common source of genotyping error resulted from varying degrees of adenylation at
459 certain markers. The most intense peak or a cluster, which the application assigns to an allele,
460 may sometimes represent adenylated fragments and sometimes non-adenylated fragments. The
461 estimated size of the same allele will therefore vary between individuals by approximately
462 one base pair. This type of variation was much more rarely induced by stuttering, the degree
463 of which is more constant.

464 In rare instances, peaks representing alleles had the same position as taller peaks in other
465 channels and were erroneously considered as resulting from crosstalk. These errors were
466 detected because neighboring peaks of similar shapes were present (indicating stuttering or
467 heterozygosity) despite the absence of peaks in other channels at their position. More
468 frequently, small artefactual peaks were not interpreted as crosstalk because their shape was
469 irregular and/or their position was slightly shifted from that of the peaks that induced
470 interference. This issue rarely affected genotyping as these peaks were generally too small to
471 be considered as alleles. Rare errors occurred in very specific situations where the length of
472 the alleles differed by only one base pair, such that the shorter peak was considered as the
473 result of adenylation. Only visual comparison with other genotypes showed that adenylation
474 was unlikely. The genotype caller does implement such check by comparing different
475 genotypes called in the same batch (S1 Text), but this check may not always be effective.
476 Finally, shorter allele dominance in heterozygotes [22], causing the peak representing the
477 longer allele to be much smaller due to a very large difference in length between alleles (> 60
478 bp), was not always properly managed. Admittedly, whether such peak should be considered
479 as an allele is difficult to determine even for experienced users.

480 Performance

481 During the evaluation, the performances of STRyper were monitored by debugging code and
482 by the profiling tools of Xcode 15 on a MacBook Pro equipped with an M1 Pro chipset and a
483 120-Hz display comprising ~6M pixels. When it came to execution speed, importing the 648
484 chromatograms took 2.45 s, i.e., 264 chromatograms were imported per second on average.
485 Application of the size standard (which involves peak assignment to DNA ladder fragment)

486 took less than 0.12 s (~5400 chromatograms per second). Allele calling of the 3600 genotypes
487 took 0.24 s (~15000 genotypes called per second). Since chromatograms/genotypes are
488 processed successively in a single execution thread, the runtime of these tasks is proportional
489 to the number of chromatograms or genotypes processed.

490 Memory usage was measured at 132 Megabytes (MB) after chromatogram import. It peaked
491 at 250 MB after selecting the 3600 called genotypes and scrolling the 3600 traces from top to
492 bottom and back. Memory usage peaked at 460 MB after selecting the 648 samples to display
493 the stacked traces at the five channels (2000 traces displayed at once, as the application does
494 not display more than 400 stacked traces per row).

495 All tasks other than those timed above were essentially instantaneous. Only the display of the
496 3600 genotypes and the 648 samples in the right pane induced a noticeable delay of about 1
497 second. Zooming and scrolling traces was generally achieved without noticeable frame drops,
498 except when zooming in/out more than about 500 traces (stacked in several rows) near their
499 full range (about 600 base pairs).

500 Discussion

501 Based on its design, features and performance, STRyper should be a valuable tool for
502 researchers who use traditional microsatellite markers. Genotyping hundreds of *Phoxinus*
503 individuals at 12 markers with STRyper proved much faster than any of my previous
504 genotyping jobs on similar data, keeping in mind that I cannot afford a comparison with
505 recent versions of commercial competing applications. This test also showed that crosstalk
506 detection and genotype calling was reasonably efficient, and could be improved upon. While
507 the underlying methods can surely be refined, I believe that substantial improvements in these
508 areas require comparisons between samples. Trained artificial intelligence has been proposed
509 for the analysis of chromatograms [23], but this approach can only be used on limited set of
510 microsatellite markers. As STRyper, nor any equivalent software, is not immune to
511 genotyping errors (reviewed in [20]) one should always visually review genotypes and
512 perform downstream corrections on exported results (e.g., [24-26]).

513 Independently of the performance of its allele caller, the main benefits of STRyper lie in its
514 streamlined user interface that is optimized for the management and inspection of hundreds of
515 chromatograms. This optimization is essential to population geneticists, who cannot spend as
516 much time on individual genotypes as forensic researchers can. Since STRyper is not

517 designed for diagnostics and must not be used for this task (it comes with no warranty), it
518 does not assume that allele calls are reviewed by several users. Therefore, it does not record
519 the history of manual corrections applied to genotypes (but still allows adding comments on
520 genotypes). Such feature would have cluttered the user interface for very little benefits for
521 most researchers.

522 Based on the reported metrics, users should not be concerned about the performance and
523 responsiveness of STRyper. The size of the database and the number of samples contained in
524 the selected folder should have little effect on the application performance and memory usage.
525 The application essentially shows tables (including its right pane), for which only the visible
526 rows, and a few others kept in cache for performance, are allocated in system memory (a
527 feature provided by the NSTableView class of the AppKit framework). Rows that are not yet
528 visible are not allocated, and those that move out view during scrolling eventually become
529 deallocated. When chromatograms are fetched using textual metadata (sample name, plate
530 well, plate name, run date, etc.), for example during a search though the whole database, only
531 that piece of data is fetched from the store and allocated in memory (a feature of the Core
532 Data framework). Fluorescence data is stored in separate objects (S1 Text) and is only fetched
533 and allocated in memory when traces are displayed.

534 As the application only uses about 460 MB when displaying 2000 traces at once – the most
535 that is allowed – memory usage should not be a concern either. The use of Apple-provided
536 frameworks (mainly AppKit, Core Data and Core Animation) contributes to the low memory
537 footprint and responsiveness of STRyper but would require a major rewrite of the GUI and
538 database-management code if the application were to be ported to non-Apple platforms.

539 However, methods related to chromatogram parsing, peak assignments (genotype calling and
540 sizing) and drawing of fluorescence curves do not heavily depend on these frameworks (they
541 mostly use functions written in plain C) and can be reused with only minor modifications.

542 From a GUI standpoint, several features of STRyper should be particularly useful to users.
543 The first is the distinction between alleles and additional peaks. Since the number of peaks
544 assigned to alleles never exceeds the marker ploidy, users should rarely need to remove peaks
545 to correct a genotype that was called, a repetitive task that proved rather tedious in my
546 previous genotyping jobs. The detection of additional peaks is optional, and these peaks can
547 be reviewed, added manually, removed, or simply ignored as they are not part of an
548 individual's genotype (they are listed and exported in a dedicated column). Theoretically,

549 additional peaks should allow genotyping polyploid species, but I have not tested STRyper for
550 this usage.

551 The second feature to underline is the implementation of fragment binning. The possibly to
552 assign off-bin peaks to alleles (bins) via drag-and-drop (**Fig 5**, left) is certainly a time saver
553 compared to typing allele names or selecting them among a long list. This task can even be
554 avoided by minimizing the offset between peak and bin locations (**Fig 3**) prior to binning, in
555 case variations in electrophoretic conditions have shifted the position of peaks relative to bins.
556 This is currently done manually by the user, but a fully automatic, or user-assisted, procedure
557 that minimizes the offset between peaks and bins (or theoretical fragment sizes) could be the
558 goal of future developments. Granted, binning can be performed automatically by
559 downstream programs like Tandem [17]. However, minimizing the offset between bins and
560 peak representing “standard” alleles should help to distinguish alleles whose size do not
561 follow the periodicity of the microsatellite repeat motive, and which may justify the creation
562 of specific bins. Tandem alerts the user about problematic alleles but does not create new
563 bins.

564 When it comes to database management, STRyper distinguishes itself by advanced search and
565 filtering capabilities, which help reviewing problematic cases, among other benefits. For
566 instance, all samples showing a particular allele at a marker can easily be retrieved across the
567 whole database and displayed. To this end, samples can be gathered in a smart folder
568 according to the name of the marker panel applied to them. Then, the list of their genotypes
569 can be filtered based on the marker name, and the allele name or size. Any new genotyped
570 sample presenting this allele would automatically appear in the smart folder.

571 Finally, the set of chromatograms contained in a folder (or a smart folder) with all its related
572 data (marker panels and bin sets, custom size standard(s), genotypes...) is easy to share, as it
573 can be transferred between instance of STRyper with a few mouse clicks and no option to set.
574 Making folder archives available alongside any publication using STRyper should help to
575 review results and to standardize the analysis of the same microsatellite markers by different
576 researchers.

577 Acknowledgements

578 I thank Dr. Douglass Hoffman from the United States National Institute of Health for his
579 advice on decoding HID files, and numerous colleagues for testing STRyper. I also thank Drs.

580 Matej Vucić and Frédéric Grandjean for giving access to *Phoxinus* chromatogram files, and
581 Dr. Romain Pigeault for reading a draft of the manuscript.

582 References

- 583 1. Hauser S. S., Athrey G., Leberg P. L. (2021) *Waste not, want not: Microsatellites*
584 *remain an economical and informative technology for conservation genetics*. *Ecol Evol.*
585 11(22):15800-14. <http://doi.org/https://doi.org/10.1002/ece3.8250>
- 586 2. Hodel R. G., Segovia-Salcedo M. C., Landis J. B., et al. (2016) *The report of my death*
587 *was an exaggeration: A review for researchers using microsatellites in the 21st century*. *Appl*
588 *Plant Sci.* 4(6). <http://doi.org/10.3732/apps.1600025>
- 589 3. Selkoe K. A., Toonen R. J. (2006) *Microsatellites for ecologists: a practical guide to*
590 *using and evaluating microsatellite markers*. *Ecol Lett.* 9(5):615-29.
591 <http://doi.org/10.1111/j.1461-0248.2006.00889.x>
- 592 4. Verbiest M., Maksimov M., Jin Y., et al. (2023) *Mutation and selection processes*
593 *regulating short tandem repeats give rise to genetic and phenotypic diversity across species*. *J*
594 *Evol Biol.* 36(2):321-36. <http://doi.org/10.1111/jeb.14106>
- 595 5. De Barba M., Miquel C., Lobréaux S., et al. (2017) *High-throughput microsatellite*
596 *genotyping in ecology: improved accuracy, efficiency, standardization and success with low-*
597 *quantity and degraded DNA*. *Mol Ecol Resour.* 17(3):492-507.
598 <http://doi.org/https://doi.org/10.1111/1755-0998.12594>
- 599 6. Barbian H. J., Connell A. J., Avitto A. N., et al. (2018) *CHIIMP: An automated high-*
600 *throughput microsatellite genotyping platform reveals greater allelic diversity in wild*
601 *chimpanzees*. *Ecol Evol.* 8(16):7946-63. <http://doi.org/https://doi.org/10.1002/ece3.4302>
- 602 7. Suez M., Behdenna A., Brouillet S., et al. (2016) *MicNeSs: genotyping microsatellite*
603 *loci from a collection of (NGS) reads*. *Mol Ecol Resour.* 16(2):524-33.
604 <http://doi.org/https://doi.org/10.1111/1755-0998.12467>
- 605 8. Liu P., Wilson P., Redquest B., et al. (2024) *Seq2Sat and SatAnalyzer toolkit: Towards*
606 *comprehensive microsatellite genotyping from sequencing data*. *Mol Ecol Resour.*
607 24(3):e13929. <http://doi.org/10.1111/1755-0998.13929>
- 608 9. Zhan L., Paterson I. G., Fraser B. A., et al. (2017) *megasat: automated inference of*
609 *microsatellite genotypes from sequence data*. *Mol Ecol Resour.* 17(2):247-56.
610 <http://doi.org/https://doi.org/10.1111/1755-0998.12561>
- 611 10. Alberto F. (2009) *MsatAllele_1.0: An R package to visualize the binning of*
612 *microsatellite alleles*. *J Hered.* 100(3):394-7. <http://doi.org/10.1093/jhered/esn110>
- 613 11. Palero F., González-Candelas F., Pascual M. (2011) *MICROSATELIGHT--pipeline to*
614 *expedite microsatellite analysis*. *J Hered.* 102(2):247-9.
615 <http://doi.org/10.1093/jhered/esq111>
- 616 12. Goor R. M., Forman Neall L., Hoffman D., Sherry S. T. (2011) *A mathematical*
617 *approach to the analysis of multiplex DNA profiles*. *Bull Math Biol.* 73(8):1909-31.
618 <http://doi.org/10.1007/s11538-010-9598-0>
- 619 13. Goor R. M., Hoffman D., Riley G. R. (2021) *Novel Method for Accurately Assessing*
620 *Pull-up Artifacts in STR Analysis*. *Forensic Science International: Genetics.* 51.
621 <http://doi.org/10.1016/j.fsigen.2020.102410>

- 622 14. Idury R. M., Cardon L. R. (1997) *A simple method for automated allele binning in*
623 *microsatellite markers*. *Genome Res.* 7(11):1104-9. <http://doi.org/10.1101/gr.7.11.1104>
- 624 15. Rosenblum B. B., Oaks F., Menchen S., Johnson B. (1997) *Improved single-strand DNA*
625 *sizing accuracy in capillary electrophoresis*. *Nucleic Acids Res.* 25(19):3925-9.
626 <http://doi.org/10.1093/nar/25.19.3925>
- 627 16. AMOS W., HOFFMAN J. I., FRODSHAM A., et al. (2007) *Automated binning of*
628 *microsatellite alleles: problems and solutions*. *Mol Ecol Notes.* 7(1):10-4.
629 <http://doi.org/https://doi.org/10.1111/j.1471-8286.2006.01560.x>
- 630 17. Matschiner M., Salzburger W. (2009) *TANDEM: integrating automated allele binning*
631 *into genetics and genomics workflows*. *Bioinformatics.* 25(15):1982-3.
632 <http://doi.org/10.1093/bioinformatics/btp303>
- 633 18. Clark J. M. (1988) *Novel non-templated nucleotide addition reactions catalyzed by*
634 *procaryotic and eucaryotic DNA polymerases*. *Nucleic Acids Res.* 16(20):9677-86.
635 <http://doi.org/10.1093/nar/16.20.9677>
- 636 19. Hauge X. Y., Litt M. (1993) *A study of the origin of 'shadow bands' seen when typing*
637 *dinucleotide repeat polymorphisms by the PCR*. *Hum Mol Genet.* 2(4):411-5.
638 <http://doi.org/10.1093/hmg/2.4.411>
- 639 20. Pompanon F., Bonin A., Bellemain E., Taberlet P. (2005) *Genotyping errors: causes,*
640 *consequences and solutions*. *Nature Reviews Genetics.* 6(11):847-59.
641 <http://doi.org/10.1038/nrg1707>
- 642 21. Vucic M., Jelic M., Klobucar G., et al. (2022) *A new set of microsatellite markers for*
643 *Phoxinus lumaireul sensu lato, Phoxinus marsilii and Phoxinus krkae for population and*
644 *molecular taxonomic studies*. *J Fish Biol.* 101(5):1225-34.
645 <http://doi.org/https://doi.org/10.1111/jfb.15194>
- 646 22. Walsh P. S., Erlich H. A., Higuchi R. (1992) *Preferential PCR amplification of alleles:*
647 *mechanisms and solutions*. *PCR Methods Appl.* 1(4):241-50.
648 <http://doi.org/10.1101/gr.1.4.241>
- 649 23. Taylor D., Powers D. (2016) *Teaching artificial intelligence to read electropherograms.*
650 *Forensic Science International: Genetics.* 25:10-8.
651 <http://doi.org/10.1016/j.fsigen.2016.07.013>
- 652 24. Wang C., Schroeder K. B., Rosenberg N. A. (2012) *A Maximum-Likelihood Method to*
653 *Correct for Allelic Dropout in Microsatellite Data with No Replicate Genotypes*. *Genetics.*
654 192(2):651-69. <http://doi.org/10.1534/genetics.112.139519>
- 655 25. De Meeûs T., Noûs C. (2022) *A simple procedure to detect, test for the presence of*
656 *stuttering, and cure stuttered data with spreadsheet programs*. *Peer Community Journal.* 2.
657 <http://doi.org/10.24072/pcjournal.165>
- 658 26. Van Oosterhout C., Hutchinson W. F., Wills D. P. M., Shipley P. (2004) *micro-checker:*
659 *software for identifying and correcting genotyping errors in microsatellite data*. *Mol Ecol*
660 *Notes.* 4(3):535-8. <http://doi.org/https://doi.org/10.1111/j.1471-8286.2004.00684.x>

661

662

663 Supporting information

664 **S1 Text. Details on peak detection and assignment, and on database management.** This
665 file describes methods for peak delineation, baseline fluorescence level subtraction,
666 determination of crosstalk, size assignment of molecular ladder fragments, detection of
667 microsatellite alleles, and an overview of the database managed by STRyper.

668

669 **S1 File. Exported results of the analysis performed to evaluate the application.** The file
670 contains a folder archive named “Phoxinus-2024.folderarchive”. This archive contains all data
671 related to the analysis of chromatograms from 324 *Phoxinus sp* individuals. The file can be
672 imported in STRyper as a folder called “Phoxinus 2024”. To do so, unzip the file if needed,
673 and import “Phoxinus-2024.folderarchive” via the “File/Import Archived Folder...” menu of
674 STRyper. See the STRyper help for more information.

675