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

2   **management**

3   Jean Peccoud

4

5   Laboratoire Écologie et Biologie des Interactions, Équipe Écologie Évolution Symbiose,

6   Université de Poitiers, UMR CNRS 7267, Bât. B31, 3 rue Jacques Fort, TSA 51106, 86073

7   Poitiers Cedex 9, France

8   jeanpeccoud@gmail.com

9
10   Running title: An application for microsatellite genotyping

11
12   **Statements relating to ethics and integrity policies**

18

19

## Abstract

21    Microsatellite markers analyzed by capillary sequencing remain useful tools for rapid

22    genotyping and low-cost studies. This contrasts with the lack of a free application to analyze

23    chromatograms for microsatellite genotyping that is not restricted to human genotyping. To

24    fill this gap, I have developed STRyper, a macOS application whose source code is published

25    under the General Public License. STRyper only uses macOS libraries, making it very

26    lightweight, responsive, and behaving like a modern application. Its three-pane window

27    enables easy management and viewing of chromatograms imported from .fsa and .hid files,

28    the creation of size standards and of microsatellite marker panels (including bins). STRyper

29    features powerful search capabilities (with smart folders) and a modern graphical user

30    interface allowing, among others, the edition of DNA ladders and of individual genotypes by

31    drag-and-drop. It also introduces a new way to mitigate the effect of variations in

32    electrophoretic conditions on estimated allele sizes.

33

34    Keywords: microsatellites, capillary electrophoresis, chromatograms, population genetics,

35    graphical user interface

## Introduction

37    More than three decades after their first use, microsatellites markers, also known as short

38    tandem repeat (STR) loci, remain popular DNA markers to assess gene flow, population

39    history, structure and membership, ancestry, or the integrity of laboratory breeding lines,

40    among other uses [1, 2]. When locus-specific variation is not the focus of a study, a limited

41    number of microsatellite markers are sufficient to assess evolutionary processes affecting

42    the whole genome and to genetically identify an individual [3]. This ability stems for the

43    sheer number of alleles per marker, which often counts in the dozens, leading to a per-locus

44    information amount that exceeds that of single-nucleotide polymorphisms (SNPs) [4].

45    Due to frequent indels in the microsatellite repeat motive, microsatellites alleles essentially

46    differ in their length, which can be estimated by simple electrophoresis of amplicons.

47    Amplicon sequencing by the Illumina technology has however emerged as relatively

48    affordable and more reliable alternative to capillary electrophoresis (e.g. De Barba, Miquel

49    (5), Barbian, Connell (6), Suez, Behdenna (7)). These pipelines of microsatellite analysis via

50    amplicon sequencing forgo the definition and optimization of microsatellite multiplexes and

51    allow the analysis of more markers. In species for which tried and tested microsatellite

52    multiplexes exist, microsatellite genotyping via electrophoresis still offers a compelling

53    money- and time-saving solution.  At a few dollars per individual in terms of consumables

54    (for a couple of multiplexes typically combining 10-20 loci) genotyping can be performed

55    locally in one day, as it amounts to DNA extraction, PCR, amplicon dilution and placing a

56    plate in a capillary sequencer. When a quick answer is needed or when only few individuals

57    need analyzing, typically for simple genotype checking, this traditional technique remains

58    the cheapest and easiest one.

59    However, the difficulty sharply rises when it comes to analyzing the results of capillary

60    electrophoresis. As opposed to genotyping via NGS, which is generally done via fully- or

61    partially automated free tools (e.g., [8, 9]), traditional microsatellite genotyping requires

62    inspecting fluorescence curves, hence applications with a complex graphical user interface

63    (GUI), which are rarely free. To various degrees, these applications are focused on human

64    identification by genotyping and forensics. A such, they are packed with features and

65    safeguards that are of little relevance to most researchers, which somewhat complicate their

66    use, and which may come at a high price.

67    This is the case of GeneMapper by ThermoFisher Scientific, a commercial application running

68    on the Windows operating system, and which remains, to my knowledge, the most widely

69    used for microsatellite genotyping. A Google scholar search for "genemapper", excluding

70    references, patents and review articles, and limited to 2023 and 2024, returned 2820 results

71    as of July 20th 2024. Most results pertained to medicine and forensics or may correspond to

72    preprints, but the first 130 results comprised 15 English-written studies on non-human

73    species using traditional microsatellite analyzes, indicating that this technique is far from

74    abandoned.

75    GeneMarker by Softgenetics is a similar commercial application. The price of a license of

76    either software may restrict its installation to a single computer per research laboratory. A

77    free alternative from ThermoFisher Scientific, Peak Scanner, has limited functionalities.

78    Complementary command-line tools [10, 11] provide missing features such as allele scoring

79    via binning, but may dissuade those who seek to conduct fragment analyses from

80    chromatogram import to the export of individual genotypes in a single user-friendly

81    application. In that regard, Geneious Prime and its microsatellite analysis plugin may

82    represent an interesting tradeoff between price and features. The cost of a subscription to

83    Geneious Prime may still appear excessive to users who do not need the features that this

84    product offers for the analysis of DNA sequences.

85    Osiris [10, 11], stands out as being a free, feature-rich and multi-platform (Windows and

86    macOS) tool for STR analysis. Yet, this software is, as far as I know, rarely used by population

87    geneticists, possibly because it is highly specialized in human identification.

88    Researchers, especially population geneticists, would therefore benefit from a free

89    application enabling quick microsatellite genotyping and management of thousands of

90    samples. To meet this need, I have developed STRyper, an open-source, lightweight and

91    user-friendly application that can analyze chromatogram files for STR genotyping. STRyper is

92    published under the GNU General Public License v. 3 and its name is a portmanteau of "STR"

93    and "Genotyper". As described below, STRyper features a modern GUI allowing, among

94    others, unconstrained chromatogram management via nested folders, advanced and

95    dynamic metadata-based chromatogram search with "smart" folders, easy folder

96    import/export, chromatogram and genotype filtering based on multiple criteria, the

97    definition of microsatellite multiplexes and custom size standards, fast and responsive

98    visualization of fluorescence curves with animated zooming and automatic vertical scaling,

99    the edition of DNA ladders and of individual genotypes by drag-and-drop, and a new way to

100   mitigate the effect of variations in electrophoretic conditions on estimated allele sizes. The

101   application and its codebase are available at https://github.com/jeanlain/STRyper.

## Methods

103   The efficiency of an application designed for microsatellite genotyping mostly relies on its

104   GUI, as chromatograms must be visually checked, and genotypes validated without a

105   command line. However, underlying methods of fluorescence data analysis for automatic

106   genotype calling are described first and in greater details than design choices, which are less

107   of a scientific mater. Those are covered in the results section, which describe the GUI.

## Fluorescence data analysis

The analysis of fluorescence data in chromatograms starts with the delineation of peaks, whose horizontal positions represent the lengths of DNA fragments. For this, a simple algorithm was developed. This algorithm (detailed in the supplementary text) determines whether a fluorescence data point (a "scan") is elevated enough, both relative to neighbor scans and in absolute fluorescence level. Peak delineation serves as a basis to subtract baseline fluorescence level, which helps peak visualization. The method developed for this task adjusts the height (fluorescence level) of a curve such that the start and end point of each peak are placed at level zero (figure S1). Although this adjustment cannot be applied on signals that are too faint to contain meaningful peaks, it has the benefit of offering two baseline subtraction modes: one that preserves absolute peak height, and one that maintains relative peak elevation compared to the baseline (Supplementary Text and Figure S1). As this method reduces background noise (Figure S1), no smoothing algorithm was implemented.

Chromatograms always contain fluorescence data from several wavelengths (channels). In multichannel fluorescence analysis, it is crucial to determine whether a peak represents a DNA fragment or interference from another channel (i.e., "crosstalk"). The method developed for this task compares the position, shape and relative size of peaks between channels, accounting for saturation of the sequencer camera (supplementary text). While certain applications alter fluorescence data to correct for pull-up due to crosstalk [11], flagging peaks resulting from crosstalk and leaving the source signal untouched was considered sufficient. These peaks are simply ignored in automatic detection of alleles and DNA ladder fragments (detailed below), although the user can manually assign these peaks, should they wish to.

## Peak assignment

Two types of peaks in chromatograms must be assigned: those that correspond to DNA ladder fragments, in the context of sample sizing, and those corresponding to alleles in the context of genotyping.

The method used to detect DNA ladder fragments and assign them to sizes of a known size standard is based on relative peak positions and accounts for non-linear relationship

138    between fragment size and migration speed (supplementary text). Peaks resulting from

139    crosstalk or whose height are unusual compared to others are ignored. To account for non-

140    linearity, a polynomial of the first, second, or third degree (depending on the user choice) is

141    used to estimate fragment size, where the response variable is the size of a fragment

142    specified in the size standard, and the explanatory variable is the scan numbers at the tip of

143    the corresponding peak (representing migration speed). This principle is also implemented in

144    other applications such as GeneMapper. Fitting is achieved via the Cholesky decomposition

145    implemented in the Linear Algebra Package (https://netlib.org/lapack/). Fitting parameters

146    are used to draw fluorescent curves ("traces") by computing the size in base pairs

147    corresponding to every scan. Traces are therefore drawn on a plot whose horizontal axis

148    represents the size in base pairs rather than the scan numbers. The horizontal distance

149    between successive scans varies unless a polynomial of the first degree (linear regression) is

150    used for the sizing.

151    To evaluate the quality of the sizing, a score from 0 to 1 was developed, based on the

152    residuals of the fitted model (differences between fragment sizes as defined in the size

153    standard, and fragment sizes estimated by the model). This score involves computing the

154    difference in residuals for every pair of adjacent peaks and is computed as follows. If $\Delta R$ is

155    the difference between residuals of every pair of adjacent peaks, $\Delta S$ the difference in scan

156    number of these peaks, $n_p$ the number of peaks and $n_s$ is number of sizes in the size

157    standard, the quality score is:

158
$$1 - \max\left(\frac{\Delta R^2}{|\Delta S|}\right)\frac{10}{3} - \frac{n_s - n_p}{10}$$

159    Any negative score is set to zero. This formula was tuned by testing many chromatograms to

160    ensure that the score is greatly reduced (often to zero) by a single assignment error, forcing

161    the user to rectify it. The score is reduced less drastically if certain sizes of the size standard

162    are not assigned to any peak, as this generally reflects problems during electrophoresis

163    rather than fixable errors.

164    Regarding allele calling, automatic genotyping must account for two main biochemical

165    processes producing DNA fragments of different lengths. One is the addition of a non-

166    template nucleotide to the 3' end of the new DNA strand by the DNA polymerase during PCR

167    [12]. Because the added nucleotide is generally an adenosine, this process is referred to as

168  "adenylation". If adenylation affects only a portion of the amplicons, they may differ in

169  length by one nucleotide, generating two peaks. The other process is "slippage" during

170  replication, causing indels in the repeated region [13]. Slippage may result in a range of

171  different amplicons whose differ by the size of the repeat, a pattern known as "stuttering".

172  These considerations served as a basis to develop a method for allele calling that first

173  identifies peak clusters resulting from these processes (detailed in the Supplementary Text),

174  and which accounts for the length of the repeat motive. In each delineated cluster, the most

175  intense peak is considered as that representing the allele. Estimation of peak intensity

176  accounts for clipping due to saturation of the fluorescence signal, in that the width of the

177  saturated region is used when peak height/area may not reflect the quantity of DNA

178  material. Stuttering and adenylation are managed internally by the application, the user

179  remains free to manually assign an allele to any peak.

180  Importantly, the method does not consider the absolute height or shape of a peak to call the

181  first allele, beyond the fact that a minimal fluorescence level is required to delineate a peak

182  (see Supplementary Text). If a peak is detected in the marker range (see below) and is not

183  interpreted as crosstalk, at least one allele will be called. It was considered that the

184  assessment of peak quality was better left to the user, who is expected to visually inspect

185  every genotype.

186  Genotyping

187  Identifying peaks representing alleles in a chromatogram requires a user-defined range of

188  expected allele sizes at a microsatellite marker. For a diploid individual, the number of

189  different alleles detected within that range determines the individual's genotype:

190  homozygous if one allele is detected, heterozygous otherwise. Because this inference is

191  invalid for polyploid markers, it was decided that only haploid and diploid markers could be

192  defined in the application, constraining the maximum number or alleles per locus to 2. To

193  cope with this constraint, the ability to annotate additional DNA fragments of interest, either

194  automatically or manually, was implemented. Additional fragments may inform on the

195  presence of paralogs, polyploidy, insufficient specificity of the PCR or contamination

196  between samples. The application therefore distinguishes two types of peaks: those that are

197  interpreted as alleles and whose number is limited to the ploidy of the marker, and others

198  representing these additional DNA fragments. Because neither should comprise fragments

199    produced by stuttering or adenylation, additional peaks are detected like alleles are (i.e., by

200    identifying peak clusters). The relative height of peaks is used to categorize alleles (higher

201    peaks) and additional peaks (smaller peaks).

202    The second crucial phase of microsatellite genotyping is the characterization of alleles based

203    on estimated amplicon sizes. Because the estimated size of an amplicon slightly varies

204    between electrophoreses and never exactly match its true length [14] amplicons are

205    assigned to alleles via "binning" [15], where bins are non-contiguous intervals delimiting the

206    expected sizes of fragments corresponding to alleles of a marker. Proper bin definition must

207    account for factors affecting amplicon mobility during electrophoresis [16], which often

208    cause the estimated distance between consecutive microsatellite alleles (in base pairs) to

209    slightly differ from the repeat motive length [15]. Binning can be left to specialized programs

210    like Tandem [17], which can work on fragment sizes estimated by other programs like

211    STRyper. The management of bins within STRyper was still considered a necessity. Indeed,

212    visualizing bins behind traces helps to characterize alleles that do not conform to the

213    periodicity of the repeat motive, and to mitigate variations in fragment sizes between

214    sequencer runs (further discussed below). Methods to import bin sets as text files (produced

215    by other tools like Tandem), to generate bin sets within the application, and to modify bins

216    individually were therefore implemented. Given a set of bins for a marker, the principle of

217    binning is simple: if the size of a fragment falls within a bin, the fragment takes the bin

218    name. By default, a bin name is the rounded size of its midpoint when the bin was created,

219    but it can be changed to any Unicode string.

220    The width of a bin might not cover the full range of estimated sizes of amplicons from given

221    allele over all electrophoretic conditions. Regularly, a peak representing an allele would fall

222    outside the corresponding bin, although identical fragments that migrated in other

223    sequencer runs were properly binned. To circumvent the issue, a mixture of amplicons of

224    known sizes for each marker, known as "allelic ladder" or "inter-lane standard", can be

225    added alongside samples for each run or sequencing plate. Allelic ladders are however only

226    available for model species.

227    To mitigate the issue, a novel approach was conceived. Rather than managing multiple bin

228    sets per marker, this approach considers that it is the sizes of alleles, not the position of bins,

229    which should be adjusted. The method thus correct fragment sizes using the formula $y = a +$

230     $bx$, where $x$ is the observed size of a DNA fragment, as estimated by the fitted model

231     mentioned earlier, $y$ is the size that will be used for fragments identified in the marker range,

232     and $a$ and $b$ are constants (hereafter called "offset parameters"). This approach assumes

233     that the effect of varying electrophoresis conditions can be compensated by this linear

234     combination. If there is no correction, $a = 0$ and $b = 1$. Good offset parameters are those that

235     minimize the distance (in base pairs) between peaks and their corresponding bins. Because

236     automatically determining which bins and peaks to associate might have been error-prone, a

237     manual GUI-based method was developed. The application lets the user move and/or resize

238     a rectangle representing the range of the bin set such that bins coincide with peaks (Figure

239     1). To infer offset parameters $a$ and $b$ from this operation, we let $s$ represents the start of a

240     bin and $e$ its end, in base pairs. If $s'$ and $e'$ represent the corresponding boundaries after the

241     user has moved the bin set appropriately, the offset parameters can be computed by solving

242
$$\begin{cases} s' = a + bs \\ e' = a + be \end{cases}$$

243     Hence $b = \dfrac{e'-s'}{e-s}$ and $a = s' - s\dfrac{e'-s'}{e-s}$ .

244     Since the user moves the bin set as a whole, the operation yields same offset parameters for

245     all bins. These parameters are then associated to the chromatograms involved in the

246     procedure (e.g., those displayed in Figure 1) and a given marker. Bin boundaries are

247     internally unchanged (they remain $e$ and $s$) such that no new bin set is created. However,

248     bins are *displayed* using $s'$ and $e'$.

## Development of the application

250     The principle described above were incorporated in a chromatogram management

251     application controlled by a GUI. The application needed to implement chromatogram file

252     parsing, importing and organization, the display and editing of metadata, performant

253     fluorescence curve drawing, the definition of microsatellite markers and bins, genotype

254     editing, the exportation of results, etc.

255     GUI development relies on application programming interfaces and frameworks that depend

256     on the target operating system and development tools. These were dictated by my use of

257     the Mac operating system (macOS) and by the fact that developing STRyper was a hobby

258     project of an evolutionary biologist, not the effort of a team of professional developers.

259  Being unencumbered by cross-platform development gave me the freedom to choose the

260  right tools to program a GUI that was intuitive, responsive and consistent with "native"

261  macOS applications. STRyper was thus developed using Xcode and frameworks provided by

262  Apple

263  (https://developer.apple.com/library/archive/documentation/MacOSX/Conceptual/OSX_Tec

264  hnology_Overview/SystemFrameworks/SystemFrameworks.html). These frameworks

265  include "Core Data", which was used to define and manage objects representing

266  chromatograms, marker panels (multiplexes), bins, alleles, genotypes and size standards,

267  and to save them in a persistent relational database (Figure S2). Internally, Core Data uses

268  the SQLite database engine to manage the persistent store. GUI elements (windows, views,

269  controls and so on) were implemented using "AppKit". "Core Graphics" functions were used

270  to draw fluorescent curves. "Core Animation" layers accelerate compositing via the graphical

271  processing unit (GPU) and provide fluid animation of the interface. These object-oriented

272  frameworks (expect Core Graphics) required the use of the Objective-C programming

273  language (a superset of C) when the project started. The application code was written in the

274  latest version (2.0) of this language.

275  ## Evaluating the application

276  STRyper was developed to facilitate the genotyping of numerous individuals (see results)

277  from chromatogram import, management and viewing, to genotyping editing and data

278  export. How well it performs at these tasks cannot be evaluated without a part of

279  subjectivity.

280  The ability of an application to assign the right peaks to a DNA ladder fragments or alleles

281  (i.e., allele calling) can be quantified more objectively by comparing these assignments to a

282  reference, which is the assignments that an experienced user would have made by visually

283  inspecting the chromatograms. Another reference, which could be used to evaluate the

284  allele caller specifically, is the genotypes obtained at the same markers from an independent

285  and more reliable method, typically amplicon sequencing. Such reference would allow

286  detecting errors that even an experienced user would not detect. These errors may arise

287  from variations in the amplification of alleles (e.g., null alleles) and motility of amplicons

288  (leading to migration speed not being proportional to fragment length), due to the intrinsic

289  properties of these fragments or variations in experimental conditions, including

290    instruments and operators (reviewed in [18]). Mitigating errors that are visually indetectable

291    should not reasonably be expected from this application. Comparing genotypes called by

292    STRyper to those obtained by sequencing would therefore not constituting a fair evaluation

293    of the allele caller, even if sequence data were available for the same individuals and

294    markers (I am not aware of the public availability of such dataset). The frequency of manual

295    corrections that an experienced user must apply to automatic peak assignment was

296    therefore used as a metric of the application performance, even though it is partly user

297    dependent.

298    Given these limitations, it was considered more valuable to evaluate STRyper as part of an

299    ongoing study (Vucić et al., in prep) instead of reanalyzing previously published data.

300    Chromatograms were obtained from 314 individuals of the freshwater fish *Phoxinus*

301    *lumaireul* (Teleostei, Cypriniformes), each amplified at two 6-plexes of microsatellite

302    markers developed by  Vucic, Jelic (19). Amplicons were submitted to electrophoresis in an

303    SeqStudio sequencer (Applied Biosystems) after addition of the GeneScan 500-LIZ size

304    standard. After importing the 648 chromatograms (628 from amplified samples and 20

305    negative controls) into STRyper, the GeneScan 500 size standard was applied to each using

306    the $3^{rd}$ degree polynomial as sizing method. Assignments of DNA ladder fragments to sizes

307    were then checked visually and corrected if necessary. For each marker, a set of bins was

308    generated in one step by a specifying a bin width of 1 base pair and setting bin spacing

309    according to the length of the microsatellite repeat motives [19]. For samples of a reference

310    sequencing plate, the bin set was moved and resized as a whole, such that bins position

311    matched peaks corresponding to alleles. In other plates, when bins of a marker and peaks

312    appeared slightly misaligned, a correction factor was applied according to the procedure

313    described in Figure 1. Individual bins were added at locations indicating the presence of

314    alleles that did not strictly follow the repeat pattern (probably due to mutations in

315    microsatellite flanking regions). Then genotypes were called, visually checked, and manually

316    corrected if necessary. The complete project was then exported as an archive

317    (supplementary file S1).

318    During this analysis, performance metrics were measured. Memory usage was monitored in

319    Xcode 15, and debugging code was added to measure the time taken by the three

320    operations that were found to cause a perceivable delay: chromatogram import, application
321    of the size standard, and allele calling.

## Results

### General characteristics of STRyper

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

328    The application comprises a main window (Figure 2) composed of three panes; a design
329    paradigm used by several database-management applications like email clients. The left
330    collapsible sidebar is a hierarchical list of folders and subfolders containing samples. Folder
331    and samples can be organized freely by drag and drop. A middle pane shows the content of
332    the selected folder (samples and associated genotypes) and comprises tabs to manage size
333    standards and markers. The right pane shows the traces (fluorescent curves) of selected
334    samples and genotypes.

335    STRyper uses very few modal panels or dialogs to validate user actions and all actions that
336    affect the database can be undone. Most are at a couple of clicks away or less as they do not
337    require opening and closing windows. Drag and drop can be used throughout: from
338    importing samples to applying size standards, markers, and to manually attributing alleles or
339    size molecular ladder fragments to peaks.

340    STRyper can import FSA files (HID file support is experimental, as the HID format
341    specifications are not public) containing data for 4 or 5 channels (fluorescent dyes). Samples
342    are imported into folders, and they can be moved or copied between folders at any time. A
343    folder and all its content, including subfolders, samples, genotypes at microsatellite markers,
344    associated marker panels (including bins) and custom size standards, can be archived and
345    transferred between instances of the application. Upon importing an archived folder, any
346    marker panel and size standard encoded in the archive is imported unless is it already in the
347    database. The imported folder therefore shows the same content as the original one.

348 Since samples are not constrained to compartmentalized projects, the application provides

349 search tools to find and gather samples from the whole database. Users can define various

350 search criteria, including run date, sizing quality, well identifier, plate name, marker panel

351 name, etc. Search results appear in "smart folders" which dynamically update their contents

352 as new samples meet the search criteria.

### Chromatogram viewing

354 Selecting a folder of the database shows all its samples, and associated genotypes if a panel

355 of microsatellite markers have been applied to the samples. Samples can be filtered and

356 sorted by various metadata items constituting columns that can be hidden and reordered.

357 An inspector panel dynamically updates to show information about selected samples,

358 including sizing information (Figure 3).

359 Upon selecting samples in the table, their chromatograms are instantaneously displayed on

360 the right pane (Figure 2). As the application fully supports the dark theme of macOS (version

361 10.14 or more recent), it can display traces on a dark background to alleviate eye strain. Any

362 region in which a peak statured the sequencer camera is shown behind curves as a rectangle

363 whose color reflects the channel that likely caused saturation. Traces can be scrolled and

364 zoomed in/out horizontally via trackpad gestures such as swipe, pinch and double tap, via

365 the scroll wheel, or by dragging the mouse over horizontal rulers to define a size range.

366 Dragging the mouse over the vertical ruler sets the fluorescence level at the top of the view,

367 hence the vertical scale. Zooming is animated, which helps users keep track of the range (in

368 base pairs) that is displayed.

369 Viewing options include automatic vertical scaling to the highest visible peaks, synchronizing

370 the vertical scales and horizontal positions, showing/hiding bins and region of fluorescence

371 saturation, stacking curves from several samples or channels in the same view, and

372 subtracting the baseline fluorescence level (see Methods). An original option fills the areas

373 under peaks resulting from crosstalk with the color of the channel that was inferred to

374 induce crosstalk. This feature helps users avoid considering these peaks as alleles or DNA

375 ladder fragments and makes clear why they were ignored during automatic genotyping.

### Applying size standards and checking molecular ladders

376 To apply size standards to samples, STRyper comes with several widely used standards,
377 namely those from the GeneScan brand. Users can easily edit these size standards within the
378 application and make their own.

380 STRyper displays the trace of the molecular ladder like any other trace, letting users switch
381 spontaneously between genotype and molecular ladder editing. Sizes attributed to
382 molecular ladder fragment can be changed by dragging and dropping size labels onto peaks.
383 Any change to the molecular ladder automatically updates the sizing of the sample without
384 user validation. The red component of the color used for size labels is proportional to the
385 difference between the computed size of a peak and its theoretical size, making size
386 assignment errors easy to spot. The sample inspector (Figure 3) also helps to find such errors
387 if points deviate from the curve representing the relationship between scan number and
388 peak size.

### Genotyping

390 Users can define their own panels of haploid or diploid microsatellite markers within
391 STRyper and organize them in folders. Markers are defined by their fluorescent dye, ploidy,
392 length of repeat motive, name, and the size range of their alleles. These attributes can be
393 changed after a marker is created, except for the first two. Markers can be copied between
394 panels. Users can export marker panels, which contain bins, to text files conforming to
395 simple specifications described in the user guide. These text files can be imported back as
396 marker panels. STRyper can also import panel and bin description text files exported from
397 Genemapper.

398 A set of automatically named bins for a marker can be added by specifying the width and
399 spacing of bins. To accommodate the fact that the observed distance between microsatellite
400 alleles slightly differs from the repeat motive [15], the position and width of the whole be
401 set can be adjusted by clicking and dragging, in a fashion similar to that described in Figure 1.
402 Individual bins can also be added and modified via click and drag. These actions do not
403 involve dedicated windows or panels, they can be performed at any time on the trace views
404 where bins are displayed (Figure 2, right).

405    All genotypes from samples of the current folder are listed in a table that can be sorted and

406    filtered according to various criteria (including allele names and sizes). This table lets users

407    quickly scan genotypes, as corresponding peaks and allele labels of the selected genotype(s)

408    appear on the right-pane. Correcting errors in allele call typically takes a single step: the user

409    can simply drag the mouse from a peak to a bin, drag an allele label from one peak to

410    another (Figure 4), or double-click a peak, which removes/attaches an allele from/to the

411    peak. Double clicking allele labels lets users enter arbitrary allele names directly above

412    peaks.

413    Genotypes and associated sample metadata can be exported as text files, or simply copied

414    from selected table rows to a text editor or a spreadsheet application.

415    Evaluation of the software

416    STRyper was used to genotype 314 *Phoxinus limaireul* individuals amplified at two

417    multiplexes of six markers each (see Methods). Ignoring electrophoresis failures that made

418    25 samples unusable, visual inspection of peak assignments to DNA ladder fragments found

419    issues in height chromatograms. In all cases, a size was not assigned to the appropriate peak

420    or the DNA ladder because the peak was missing or abnormally short. As issues due to

421    missing peaks cannot be fixed, manual corrections were applied to only four

422    chromatograms. Overall, the verification of the DNA ladder for all chromatograms took less

423    than five minutes.

424    For certain sequencing plates and markers (PHOX4, PHOX11, PHOX29, PHOX33 and CtoA-247

425    [19]), peaks and bins appeared slightly misaligned (by less than 0.5 base pairs). Offset

426    parameters for peak sizes were thus defined according to the procedure described in Figure

427    1. This procedure made bins coincide neatly with peak locations for all regularly spaced

428    alleles. I therefore saw no evidence that the use of linear relationship to estimate offset

429    parameters was inappropriate.

430    In cases of PCR failures, the application assigned relatively faint peaks amounting to noise as

431    alleles, which was expected. These cases were easily detected by visual inspection.

432    The only common source of error resulted from varying degrees of adenylation at certain

433    markers. The most intense peak or a cluster, which the application assigns to an allele, may

434    sometime represent adenylated fragments and sometimes non-adenylated fragments. The

435    estimated size of the same allele will therefore vary between individuals by approximately
436    one base pair. This type of variation was much more rarely induced by stuttering, the degree
437    of which is more constant.

438    In rare instances, peaks representing alleles had the same position as taller peaks in other
439    channels and were erroneously considered as resulting from crosstalk. These errors were
440    detected because neighboring peaks of similar shapes were present (indicating stuttering or
441    heterozygosity) despite the absence of peaks in other channels at their position. More
442    frequently, small artefactual peaks were not interpreted as crosstalk because their shape
443    was irregular and/or their position was slightly shifted from that of the peaks that induced
444    interference. This issue rarely affected genotyping as these peaks were generally too small
445    to be considered as alleles. In a few cases though, an artifactual peak (not identified as such)
446    was taken as an allele instead of the correct peak. This occurred in very specific situations
447    where the length of the true alleles differed by only one base pair, such that one of the
448    peaks they induced was considered to result from adenylation. Only comparison with other
449    genotypes showed that adenylation was unlikely. The genotype caller does implement such
450    check (Supplementary Text) and tries to correct genotypes that were initially considered
451    homozygous. In the cases described here however, the genotypes were considered
452    heterozygous, since an artefactual peak was taken as the second allele, and were therefore
453    not checked.

454    Finally shorter allele dominance in heterozygotes [20], causing the peak representing the
455    longer allele to be much smaller due to a very large difference in length between alleles (>
456    60 bp), was not always properly managed.  Admittedly, whether such peak should be
457    considered as an allele is difficult to determine even for experienced users.

458    Speed, memory usage and responsiveness

459    When it came to execution speed, importing the 648 chromatograms took 2.45 s, i.e., 264
460    chromatograms were imported per second on average. Application of the size standard
461    (which involves peak assignment to DNA ladder fragment) took less than 0.12 s (~5400
462    chromatograms per second). Allele calling of the 3600 genotypes took 0.24 s (~15000
463    genotypes called per second). Since chromatograms/genotypes are processed successively in
464    a single execution thread, the runtime of these tasks is proportional to the number of
465    chromatograms or genotypes processed.

466   Memory usage was measured at 132 Megabytes (MB) after chromatogram import. It peaked

467   at 250 MB after selecting the 3600 called genotypes and scrolling the 3600 traces from top

468   to bottom and back. Memory usage peaked at 460 MB after selecting the 648 samples to

469   display the stacked traces at the five channels (2000 traces displayed at once, as the

470   application does not display more than 400 stacked traces per row).

471   All tasks other than those timed above were essentially instantaneous. Only the display of

472   the 3600 genotypes and the 648 samples in the right pane induced a noticeable delay of

473   about 1 second. On a laptop with a high-resolution 120 Hz display and an M1 Pro chipset,

474   zooming and scrolling traces was generally achieved without noticeable frame drops, except

475   when zooming in/out more than about 500 traces (stacked in several rows) near their full

476   range (about 600 base pairs).

## Discussion

478   Based on its design, features and performance, STRyper should be a valuable tool for

479   researchers who use traditional microsatellite markers. Genotyping hundreds of *Phoxinus*

480   individuals at 12 markers with STRyper proved much faster than any of my previous

481   genotyping jobs on similar data, keeping in mind that I cannot afford a comparison with

482   recent versions of commercial competing applications. This test also showed that crosstalk

483   detection and genotype calling was reasonably efficient, although improvable. While the

484   underlying methods can surely be refined, I believe that substantial improvements in these

485   areas require comparisons between samples. Ideally, trained artificial intelligence could be

486   employed to analyze chromatograms [21], assuming STRyper development benefits from the

487   contribution of AI specialists. As STRyper is not immune to genotyping errors (reviewed in

488   [18]) one should always visually review genotypes and perform downstream corrections on

489   exported results (e.g., [22-24]).

490   Independently of the performance of its allele caller, the main benefits of STRyper lie in its

491   streamlined user interface that is optimized for the management and inspection of hundreds

492   of chromatograms. This optimization is essential to population geneticists, who cannot

493   spend as much time on individual genotypes as forensic researchers can. Since STRyper is

494   not designed for diagnostics and must not be used for this task (it comes with no guaranty),

495   it does not assume that allele calls are reviewed by several users. Therefore, it does not

496    record the history of manual corrections applied to genotypes (but still allows adding

497    comments on genotypes). Such feature would have cluttered the user interface for very little

498    benefits for most researchers.

499    Based on the reported metrics, users should not be concerned about the performance and

500    responsiveness of STRyper. The size of the database and the number of samples contained in

501    the selected folder should have little effect on the application performance and memory

502    usage. The application essentially shows tables (including its right pane), for which only the

503    visible rows, and a few others kept in cache for performance, are allocated in system

504    memory (a feature provided by the NSTableView class of the AppKit framework). Rows that

505    are not yet visible are not allocated, and those that move out view during scrolling

506    eventually become deallocated. When chromatograms are fetched using textual metadata

507    (sample name, plate well, plate name, run date, etc.), for example during a search though

508    the whole database, only that piece of data is fetched from the store and allocated in

509    memory (a feature of the Core Data framework). Fluorescence data is stored in separate

510    objects (Supplementary Text) and is only fetched and allocated in memory when traces are

511    displayed.

512    As the application only uses about 460 MB when displaying 2000 traces at once – the most

513    that is allowed – memory usage should not be a concern either. The use of Apple-provided

514    frameworks (mainly AppKit, Core Data and Core Animation) contributes to the low memory

515    footprint and responsiveness of STRyper but would require a major rewrite of the GUI and

516    database-management code if the application were to be ported to non-Apple platforms.

517    However, methods related to chromatogram parsing, peak assignments (genotype calling

518    and sizing) and drawing of fluorescence curves do not heavily depend on these frameworks

519    (they mostly use functions written in plain C) and can be reused with only minor

520    modifications.

521    From a GUI standpoint, several features of STRyper should be particularly useful to users.

522    The first is the distinction between alleles and additional peaks. Since the number of peaks

523    assigned to alleles never exceeds the marker ploidy, users should rarely need to remove

524    peaks to correct a genotype that was called, a repetitive task that proved rather tedious in

525    my previous genotyping jobs. The detection of additional peaks is optional, and these peaks

526    can be reviewed, added manually, removed, or simply ignored as they are not part of an

527  individual's genotype (they are listed and exported in a dedicated column). Theoretically,

528  additional peaks should allow genotyping polyploid species, but I have not tested STRyper

529  for this usage.

530  The second feature to underline is the implementation of fragment binning. The possibly to

531  assign off-bin peaks to alleles (bins) via drag-and-drop (Figure 4, left) is certainly a time saver

532  compared to typing allele names or selecting them among a long list. This task can even be

533  avoided by minimizing the offset between peak and bin locations (Figure 1) prior to binning,

534  in case variations in electrophoretic conditions have shifted the position of peaks relative to

535  bins. This is currently done manually by the user, but a fully automatic, or user-assisted,

536  procedure that minimizes the offset between peaks and bins (or theoretical fragment sizes)

537  could be the goal of future developments. Granted, binning can be performed automatically

538  by downstream programs like Tandem [17]. However, minimizing the offset between bins

539  and peak representing "standard" alleles should help to distinguish alleles whose size do not

540  follow the periodicity of the microsatellite repeat motive, and which may justify the creation

541  of specific bins. Tandem alerts the user about problematic alleles but does not create new

542  bins.

543  When it comes to database management, STRyper distinguishes itself by advanced search

544  and filtering capabilities, which help reviewing problematic cases, among other benefits. For

545  instance, all samples showing a particular allele at a marker can easily be retrieved across

546  the whole database and displayed. To this end, samples can be gathered in a smart folder

547  according to the name of the marker panel applied to them. Then, the list of their genotypes

548  can be filtered based on the marker name, and the allele name or size. Any new genotyped

549  sample presenting this allele would automatically appear in the smart folder.

550  Finally, the set of chromatograms contained in a folder (or a smart folder) with all its related

551  data (marker panels and bin sets, custom size standard(s), genotypes…) is easy to share, as it

552  can be transferred between instance of STRyper with a few mouse clicks and no option to

553  set. Making folder archives available alongside any publication using STRyper should help to

554  review results and to standardize the analysis of the same microsatellite markers by

555  different researchers.

# Acknowledgements

# References

1. Hauser SS, Athrey G, Leberg PL. *Ecol. Evol.* 2021;11(22):15800-14. http://doi.org/https://doi.org/10.1002/ece3.8250

2. Hodel RG, Segovia-Salcedo MC, Landis JB, Crowl AA, Sun M, Liu X, et al. *Appl Plant Sci*. 2016;4(6). http://doi.org/10.3732/apps.1600025

3. Selkoe KA, Toonen RJ. *Ecol Lett*. 2006;9(5):615-29. http://doi.org/10.1111/j.1461-0248.2006.00889.x

4. Verbiest M, Maksimov M, Jin Y, Anisimova M, Gymrek M, Bilgin Sonay T. *J Evol Biol*. 2023;36(2):321-36. http://doi.org/10.1111/jeb.14106

5. De Barba M, Miquel C, Lobréaux S, Quenette PY, Swenson JE, Taberlet P. *Mol. Ecol. Resour.* 2017;17(3):492-507. http://doi.org/https://doi.org/10.1111/1755-0998.12594

6. Barbian HJ, Connell AJ, Avitto AN, Russell RM, Smith AG, Gundlapally MS, et al. *Ecol. Evol.* 2018;8(16):7946-63. http://doi.org/https://doi.org/10.1002/ece3.4302

7. Suez M, Behdenna A, Brouillet S, Graça P, Higuet D, Achaz G. *Mol. Ecol. Resour.* 2016;16(2):524-33. http://doi.org/https://doi.org/10.1111/1755-0998.12467

8. Liu P, Wilson P, Redquest B, Keobouasone S, Manseau M. *Mol Ecol Resour*. 2024;24(3):e13929. http://doi.org/10.1111/1755-0998.13929

9. Zhan L, Paterson IG, Fraser BA, Watson B, Bradbury IR, Nadukkalam Ravindran P, et al. *Mol. Ecol. Resour.* 2017;17(2):247-56. http://doi.org/https://doi.org/10.1111/1755-0998.12561

10. Goor RM, Forman Neall L, Hoffman D, Sherry ST. *Bull Math Biol*. 2011;73(8):1909-31. http://doi.org/10.1007/s11538-010-9598-0

585    11.    Goor RM, Hoffman D, Riley GR. *Forensic Science International: Genetics*. 2021;51.
586    http://doi.org/10.1016/j.fsigen.2020.102410

587    12.    Clark JM. *Nucleic Acids Res*. 1988;16(20):9677-86.
588    http://doi.org/10.1093/nar/16.20.9677

589    13.    Hauge XY, Litt M. *Hum Mol Genet*. 1993;2(4):411-5.
590    http://doi.org/10.1093/hmg/2.4.411

591    14.    Rosenblum BB, Oaks F, Menchen S, Johnson B. *Nucleic Acids Res*. 1997;25(19):3925-9.
592    http://doi.org/10.1093/nar/25.19.3925

593    15.    Idury RM, Cardon LR. *Genome Res*. 1997;7(11):1104-9.
594    http://doi.org/10.1101/gr.7.11.1104

595    16.    AMOS W, HOFFMAN JI, FRODSHAM A, ZHANG L, BEST S, HILL AVS. *Mol. Ecol. Notes*.
596    2007;7(1):10-4. http://doi.org/https://doi.org/10.1111/j.1471-8286.2006.01560.x

597    17.    Matschiner M, Salzburger W. *Bioinformatics*. 2009;25(15):1982-3.
598    http://doi.org/10.1093/bioinformatics/btp303

599    18.    Pompanon F, Bonin A, Bellemain E, Taberlet P. *Nature Reviews Genetics*.
600    2005;6(11):847-59. http://doi.org/10.1038/nrg1707

601    19.    Vucic M, Jelic M, Klobucar G, Jelic D, Gan HM, Austin C, et al. *J. Fish Biol.*
602    2022;101(5):1225-34. http://doi.org/https://doi.org/10.1111/jfb.15194

603    20.    Walsh PS, Erlich HA, Higuchi R. *PCR Methods Appl*. 1992;1(4):241-50.
604    http://doi.org/10.1101/gr.1.4.241

605    21.    Taylor D, Powers D. *Forensic Science International: Genetics*. 2016;25:10-8.
606    http://doi.org/10.1016/j.fsigen.2016.07.013

607    22.    Wang C, Schroeder KB, Rosenberg NA. *Genetics*. 2012;192(2):651-69.
608    http://doi.org/10.1534/genetics.112.139519

609    23.    De Meeûs T, Noûs C. *Peer Community Journal*. 2022;2.
610    http://doi.org/10.24072/pcjournal.165

611    24.    Van Oosterhout C, Hutchinson WF, Wills DPM, Shipley P. *Mol. Ecol. Notes*.
612    2004;4(3):535-8. http://doi.org/https://doi.org/10.1111/j.1471-8286.2004.00684.x

613

## Data accessibility

STRyper and its source code are available at https://github.com/jeanlain/STRyper/

The archive of the folder containing the analyzed data from the 314 *Phoxinus sp* individuals is available as supplementary file S1.

## Author Contributions

JP developed the application, analyzed data and wrote the paper.

621    Figures



622

623    Figure 1. A case of out-of-bin alleles that is solved. Both images show the stacked traces
624    from 8 samples of the same sequencer run. Peaks represent amplicons of a dinucleotide
625    marker called "Ap22". Its range is represented by a horizontal red segment above the ruler
626    showing graduations in base pairs (bp). Bins appear as grey rectangles. Top: peaks are
627    shifted to the left with respect to bins, and more so for longer alleles, although bins are
628    separated by exactly two base pairs. Bottom: the user has moved and narrowed the light-
629    pink rectangle representing the range of the marker, such that bins coincide with peaks. This
630    move translates into offset parameters $a$ = −6.40 and $b$ = 1.029 (see main text). As a result,
631    the estimated size of peaks overlapping bin 258 (bottom image) has changed from ~257 bp
632    to ~258.1 bp.

633

634

Figure 2. The main window of STRyper. The left pane contains the list of folders and smart folders (search results) containing samples. The middle pane is a split view comprising a top pane listing the samples of the selected folder. Its bottom pane has four tabs, which are from left to right: an inspector showing data on selected samples (Figure 2), a table of genotypes from the samples shown on the top pane, the marker library (currently shown) and the size standard library. The right pane shows the traces of selected samples in a scrollable view that can display thousands of traces. The red channel currently shows the range of a diploid DNA marker ("PHOX29") that contain bins shown as vertical grey rectangles. Alleles are annotated with rectangular labels colored after the channel. Supplementary peaks in bins "267" and "259" have been annotated. The orange channel shows the molecular ladder.

646

Figure 3. The sample inspector of STRyper. This panel with three collapsible sections
dynamically updates to display information on samples that are selected in the sample table
(Figure 2). The plot at the bottom shows the relationship between the time at which DNA
fragments of the molecular ladder (black crosses) were detected by the sequencer camera
(the X axis) and their observed sizes in base pairs (the Y axis). The relationship used to
estimate fragment sizes is established by fitting a polynomial (here, of the third degree) to
the points shown on the plot. This polynomial is represented by the orange curve.

654

655



656

Figure 4. Genotype editing by drag and drop in STRyper. Vertical grey rectangles represent
bins that define expected ranges of microsatellite alleles. Each bin has a name displayed on

659    top. Allele names are represented by colored labels above peaks. Left-hand screen capture:

660    the user is dragging the mouse from a peak to a bin. This will assign the peak an allele named

661    after the bin, thereby replacing the question mark used for alleles that are out of bins.

662    During the operation, a grey-colored handle connects the mouse location to another point

663    horizontally located at the peak tip and vertically located at the clicked point. Right-hand

664    screen capture: the user has decided that only the peak on the left should represent an

665    allele and is dragging an allele label from the right-hand peak to the other. These actions are

666    assisted by "magnetism" to lock the handle or allele label to the closest suitable destination,

667    which triggers haptic feedback on the trackpad.

## 3    Peak delineation

4    To delineate peaks in the fluorescence data, STRyper uses a simple method that enumerates

5    fluorescence levels from the first to the last recorded scan. A scan is a data point that is

6    denoted by an integer index varying from 0 to the total number of data points.

7    The method records the lowest fluorescence level ($l$) and the highest level ($h$), and their

8    respective scan numbers ($s_l$, $s_h$), observed up to the current scan number ($s_f$) whose

9    fluorescence level is denoted as $f$. A peak is delineated if $h > t$, $l/h \leq r$ and $f/h \leq r$, $t$ being the

10   minimal fluorescence level to consider a peak (by default, 100 fluorescence units) and $r$

11   being a parameter denoting the minimum peak elevation above the background.

12   Horizontally, the peak starts at scan $s_l$, and its tip is at scan $s_h$. Its right boundary will

13   correspond to the left boundary of the next peak. This method thus generates contiguous

14   peaks.

15   For best results, it was found that three rounds of peak detection should be applied to the

16   data, each round being followed by one pass of baseline fluorescence level subtraction (see

17   next section). The first two rounds use a value of 0.7 for $r$, a modest peak elevation that

18   allows the detection of faint peaks. The last iteration uses $r = 0.5$, which means that a peak

19   must be at least twice higher than the background level, considering that baseline

20   fluorescence level subtraction makes peak stand-out more.

21   After these three rounds, the left and right boundaries of each peak are delineated by the

22   closest scan from each side of the peak's tip that has a fluorescence level of 0, using

23   fluorescence levels with baseline level subtracted. This produces non-contiguous peaks.

## 24   Baseline fluorescence level subtraction

25   STRyper subtracts the baseline fluorescence level of a trace after peaks are delineated (see

26   previous section) as follows. A virtual line segment is drawn from the start to the end of each

27   peak (Figure S1). For a given scan number $s_f$, the height of the segment is denoted as $y$ and is

28   considered the "baseline fluorescence". The recorded fluorescence level for the scan is

29   denoted as $f$ and the fluorescence level at the peak tip is denoted as $h$.

30   For each value of $s_f$ within the peak, a value $v$ to subtract to the fluorescence level depends

31   on whether the absolute height of peaks should be preserved. If so, $v = y(h - f)/(h - y)$.

32   Otherwise, $v = y$. If $v$ is negative, it is set to 0. The new value for the fluorescence level is $f -$

33   $v$. After this operation, each peak starts and ends at a fluorescence level of 0.

34   The same operation is performed between peaks to reduce the background noise. Between

35   peaks, $v = y$.



36

37   **Figure S1.** Subtraction of baseline fluorescence level. A) Principle of the method. Symbols are

38   defined in the supplementary text. B) Effect of the three passes of the method (see

39   supplementary text) on fluorescence curves. Top: raw fluorescence data. Bottom:

40   fluorescence data after baseline fluorescence level was subtracted.

# Determination of crosstalk

42   STRyper determines whether a peak in fluorescence results from interference between

43   channels, i.e., crosstalk. This inference relies on the presence of saturation, or of higher peak

44   of similar shapes, in other channels.

45   A chromatogram file lists each scan number for which the signal saturated the sequencer

46   camera but does not specify which channel caused the saturation. STRyper determines this

47   channel by first delineating regions composed of consecutive scan numbers where

48   saturation occurred.

49  For each region, the channel that is considered to have caused saturation is the one whose

50  fluorescence level is the highest at the first scan of the region. This criterion does not

51  compare maximum/average fluorescence levels over the region between channels, because

52  the peak at the channel that caused saturation if often clipped and may be smaller than

53  peaks of other channels in the region. However, this peak has the highest fluorescence level

54  at the point where saturation began.

55  A peak is considered to results from crosstalk if the following conditions are met: (i) its tip

56  lies within a region where saturation is caused by another channel, and (ii) the fluorescence

57  level at the peak tip is at least twice those recorded at the scan preceding the start and the

58  scan after the end of the region. Criterion (ii) accounts for the fact that several DNA

59  fragment may have migrated at the same speed, such that legit peaks appear at the same

60  locations. However, the fluorescence level at a peak resulting from crosstalk should not be

61  high before the saturation from another channel is recorded.

62  Alternatively, crosstalk may cause a "crater" in other channels, that is, sharp peaks at the

63  edges of the saturated region. If a small peak lies near such edge and sharply decreases

64  within the saturated region, the peak is considered to result from crosstalk.

65  If a focus peak is not considered are resulting from crosstalk based on these criteria, the

66  program inspects other channels to find the one with highest fluorescence level at the peak

67  tip, and for which the fluorescence level is at least 1.6 times that at the peak tip. If it finds

68  one, it then evaluates how much peaks of both channels overlap, using two criteria. The

69  program first scales down the higher peak such that is elevation corresponds to the smaller.

70  It then measures the peak areas by summing fluorescence levels. The first criterion is

71  considered passed if the area representing the intersection between peaks is at least 30% of

72  the area representing the union of the peaks. The second criterion precisely evaluates how

73  much the peak horizontal positions are aligned. For that, the difference in fluorescence level

74  (curve height) between channels is computed at each scan along the range encompassing

75  both peaks. The sign of the difference is reversed if the scan is greater than the scan of a

76  given peak's tip. For each peak, these differences are summed across all scans of the range.

77  The second criterion is considered passed if the absolute value of each sum is less than 30%

78  the combined areas of the peaks. If both criteria are met, the program checks if other peaks

79  in the channel that may have induced crosstalk also induced crosstalk in the focus channel.

3

80  This inspection relies on the expected ratio of peak heights between the two channels,

81  which should be rather constant in the case of crosstalk and in the absence of saturation. If

82  another peak does not appear to have induced crosstalk, then the peak under consideration

83  is not considered to result from crosstalk.

## Size assignment of molecular ladder fragments

85  The algorithm conceived to assign sizes to molecular ladder fragments inspects peak in the

86  appropriate channel, ignoring those resulting from crosstalk (see previous section). In the

87  following, the "scan number" of a peak refers to the scan at its tip.

88  Peaks are first enumerated by decreasing scan numbers, and the average peak height is

89  computed at each step. Any peaks whose height is at least twice the current average and

90  whose scan number is less than 1/3 total number of scans in the trace is discarded. This

91  eliminates high-intensity peaks of short size (in base pairs) resulting from degradation of the

92  molecular ladder.

93  The algorithm then discards weak peaks amounting to "noise", which sometimes affect the

94  data. To do so, remaining peaks are enumerated by decreasing height. The enumeration

95  stops when the number of enumerated peaks corresponds to the number of sizes specified

96  in the size standard, or when a peak is at least three times smaller than the previous one.

97  Any peak that is at least twice as small as the least enumerated peak is discarded.

98  To assign remaining peaks to sizes defined in the size standard, peaks are ordered by

99  increasing scan number. The method assigns the lowest size to the first peak, and the largest

100 size to the last peak. To understand the process, picture a straight line of equation $y = a + bx$

101 passing through these two peaks on a plot where the x axis represents scan numbers, and

102 the y axis sizes in base pairs.

103 Peaks are then enumerated in decreasing order, starting from the second-to-last. The size of

104 the fragment causing a peak is estimated as $a + bx$, $x$ being the peak scan number. The size

105 of the size standard that is the closest to the observed size is assigned to the peak, only if the

106 difference between both sizes is less than 15 bp in absolute value.

107 The next peak is evaluated in the same fashion. If it is assigned to the same size as a previous

108 peak, both peaks are confronted to retain the one whose predicted size is the closest. The $a$

109 and $b$ parameter are updated to correspond to the line connecting the two peaks that were

110 assigned last. Hence, the size/scan relationship dynamically changes to account for non-

111 linearity.

112 At the end of the procedure, the shortest size of the size standard may be assigned to a

113 different peak than the one of lowest scan number. This is not the case for the longest size,

114 which remains assigned to the peak of largest scan number, although this assignment might

115 be erroneous (this is addressed using subsequent iterations, as described below).

116 A quality index is computed to evaluate the assignments. This index relies on the residuals

117 of the linear regression between scan number and size in base pairs, using ordinary least

118 squares. For each pair of successive points (peaks), the difference between residuals is

119 divided by the difference between scan numbers, both in absolute value. The mean of these

120 ratios is computed. The inverse of this mean, multiplied by the percentage of sizes that were

121 assigned to peaks, constitutes the quality index. If this index is higher than a certain value

122 (chosen at 100), the number of assigned sizes is recorded as a reference.

123 Further iterations of assignments are performed by decrementing the longest assignable size

124 (to consider the possibility that electrophoresis failed or stopped before the last fragment

125 was detected), then by decrementing the last assignable peak. Assignments are not recorded

126 if the number of assigned sizes is lower than the reference, and iterations stop when the

127 number of assignable sizes/peaks is lower than the reference.

128 In the end, the set of assignments that yielded the best quality index is retained.

## Detection of microsatellite alleles

130 To identify microsatellite alleles at a marker in a chromatogram, peaks found in the marker

131 range are first sorted by decreasing number of saturated scans they induced, then by

132 decreasing height (fluorescence level). This sorting accounts for clipping due to saturation of

133 the fluorescence signal. Hereafter, a peak position/size (in base pairs) refers to the position

134 of its tip, hence the estimated length of the DNA fragment that induced the peak.

135 For each peak (hereafter called a "reference" peak), neighboring peaks are successively

136 inspected at increasing distance to identify peak clusters. Neighbors that are at the left

137 (lower scan numbers) are inspected before those at the right. Briefly, the inspection first

138     evaluates if a neighbor resulted from stuttering: the distance between the neighbor and the

139     reference peak (of from a previously inspected neighbor already considered as a stutter)

140     must be the motive length of the marker ± 0.5 bp. In addition, the neighbor must be smaller

141     than the reference peak or than a previously inspected neighbor already considered as a

142     stutter. If these requirements are met, the neighbor is flagged as a "child" of the reference

143     peak, that is, both are considered part of the same cluster and have arisen from

144     amplification of the same allele. If the distance between a neighbor and the reference peak

145     (or a previously inspected neighbor) lies between 0.5 and 1.5 bp, and if the neighbor is

146     smaller than the reference peak, the neighbor is interpreted as resulting from adenylation of

147     the amplicon, hence as a child peak. This is also the case if the neighbor is distant from the

148     motive length ± 0.5 bp from a previous peak considered as resulting from adenylation.

149     The application stops the inspection of neighbors (at the left or at the right or a reference

150     peak) if a neighbor is already flagged as a child from a reference peak inspected prior, or if its

151     distance from the last inspected neighbor exceeds the motive length + 0.5 pb, in which case

152     it the neighbor is not considered as a child.

153     Additional checks based on peak heights are implemented to avoid considering alleles

154     differing by only one repeat motive as part of the same cluster. One check accounts for short

155     allele dominance in heterozygotes: the fact that the longer allele is almost always amplified

156     with lower yield during PCR. The resulting peak is therefore smaller, but it must not be

157     considered as a result a stuttering. To account for dropout, a neighboring peak at the right of

158     the reference peak is considered as resulting from stuttering only if its height is <30% that of

159     the reference peak. Conversely, the method also considers rare cases where the shorter

160     allele has amplified with lower yield than the longer. If allele lengths differ by just one repeat

161     motive, the left peak may be erroneously considered as a stutter. To avoid this, the method

162     computes the ratio of peak heights (left peak / right peak). If the ratio is ≥ 0.7, the program

163     looks for an additional stutter peak at the left, at distance that is the motive length ± 0.5.

164     This check assumes that a first peak arising from stuttering is followed by others with similar

165     height ratio between neighboring peaks. If no such peak is found, the left peak is not

166     considered as a child peak, and the inspection of neighbors that are at the left of the

167     reference peak stops. Note that if the shorter allele amplifies with a much lower yield than

168    the longer allele, distinguishing it from a stutter would require comparing individuals, which

169    the application does not.

170    After neighbors are inspected for all reference peaks, all peaks are processed by decreasing

171    height. The first peak is always considered as an allele and any subsequent peak will be as

172    well if the following conditions are fulfilled. First, the number of alleles must not exceed the

173    marker ploidy and the peak must not be a child peak. Then, a ratio is computed by dividing

174    its height with that of the last peak considered as an allele. If the focus peak is at the left of

175    the allele and the ratio is ≥ 0.7, it is considered as another allele. If the subsequent peak is at

176    the right, the ratio must exceed 0.3, or the ratio multiplied by the absolute difference in

177    peak positions (in bp) must exceed 4. This second condition allows alleles that are much

178    longer to yield peaks that are much smaller, accounting for short allele dominance. A peak

179    not considered as an allele is flagged as a "supplementary" peak (see main text) if it is not a

180    child peak and if its height is ≥20% of the height of last peak considered as an allele (or ≥12%

181    of the allele height if the focus peak itself has child peaks).

182    The method also considers a special case where two alleles of a heterozygote differ by only

183    one nucleotide in length. In this case, one allele may be wrongly considered as a peak

184    resulting from adenylation and a homozygous genotype would be called. Only comparison

185    with other genotypes may indicate whether the degree of adenylation of the marker is

186    compatible with this assessment. Therefore, during alle call, the application records the ratio

187    between the height of a given peak and the height of a peak inferred as an allele, if the peak

188    is within 1.5 bp from the allele. This ratio is called $l$ if the fragment represented by the peak

189    is longer than the allele or $s$ if the fragment is shorter. $L = \max(l)$ and $S = \max(s)$ are

190    computed for a given genotype. In the absence of detectable adenylation, these maxima

191    would be zero. The peak that has the highest ratio is considered as a "possible allele" of the

192    genotype.

193    If, during the same allele call, at least two genotypes were considered as heterozygous and

194    at least one as homozygous for a given diploid marker, the application computes arithmetic

195    means $M(L)$ and $M(S)$ over all genotypes that are heterozygous. The application then

196    inspects the possible allele of every homozygous genotype. If a possible allele fulfills

197    $(l > 6 \times M(L)$ and $l \geq 0.5)$ or $(s > 6 \times M(S)$ and $s \geq 0.5)$, it is promoted as allele.

## Overview of the database managed by STRyper

STRyper manages a database composed of objects of different classes that are represented in Figure S2.

An object of class *Marker* describes a microsatellite marker by specifying its *name*, and the *start* and *end* of the range of expected alleles expressed in base pairs (these attributes are inherited from a superclass called *Region*). *motiveLength* specifies the length of the repeat motive, *channel* the channel (wavelength) used to reveal amplicons (blue, green, black/yellow, or red) and *ploidy* is self-explanatory. The set called *bins* points to the bins defined for the marker. The *start* and *end* attributes of a Bin (inherited from the Region class) specify the expected range of an allele. A marker also points to the *Panel* (i.e., the multiplex) it belongs to, which reciprocally points to its markers via a set called *markers*.

An object of the *Chromatogram* class stores data imported from a .fsa or .hid file into various attributes, including the sample name (*sampleName* attribute), the *plate* name, the *well* identifier, the time of the end of sequencer run (*runStopTime*), the number of scans recorded (*nScans*), the indices of scans for which the camera was saturated (*offScaleScans*), etc. The *panel* relationship points to the panel of markers that were amplified to generate the chromatogram. Reciprocally, a panel points to all chromatograms that use it for genotyping, via a set called *samples*.

Via the *traces* set, a chromatogram points to four or five objects of class *Trace*, each of which encodes the raw fluorescence data (in the *rawData* attribute) measured at a *channel*. Peaks identified in each trace are stored in the *peaks* attribute, which is an array of structures composed of three integers: the scan at the start of the peak, the number of scans from the start to the tip, and the number of scans from the tip to the end of the peak. The *dyeName* attribute stores the name of dye that emitted the recorded fluorescence used (e.g., "6-FAM", "LIZ") and *isLadder* tells whether the trace represents the molecular ladder that was added to the sample before electrophoresis.

The size standard defining the molecular ladder is referred to by the chromatogram, via its *sizeStandard* relationship. A SizeStandard object has a *name* (e.g., "GeneScan 500") and a set of *SizeStandardSize* objects (called *sizes*), each of which has a *size* attribute specifying a

227     fragment size in base pairs. This set facilitates adding, removing or changing sizes. The

228     *editable* attribute of a size standard tells whether its *sizes* can be edited by the user.

229     The *fragments* relationship of a trace points to *LadderFragment* objects. Such object defines

230     a DNA fragment that produced a peak identified in the trace. Its *scan* attribute refers to the

231     scan at the tip of the peak (hence the location of the fragment in the trace), and its *size*

232     attribute is the size (in base pairs, taken from the size standard) that was attributed by the

233     method described in section "Size assignment of molecular ladder fragments". The *name* of

234     the fragment (shown to the user) is its *size* encoded as characters. The *offset* attribute is the

235     difference (in base pairs) between the location of the *scan* (computed via the *coefs* attribute

236     of the chromatogram, see below) and the *size*. This helps users detect size assignment

237     errors. The coefficients of the polynomial that was fitted based on detected ladder

238     fragments (see main text section "**Erreur ! Source du renvoi introuvable.**") are stored as an

239     array of floating-point numbers in the *coefs* attribute of the chromatogram to which the

240     trace belongs.

241     For a trace that is not a molecular ladder, the *fragments* set can only contain *Allele* objects.

242     An Allele defines a DNA fragment that produced a peak in the range of a marker (from the

243     panel applied to the trace's chromatogram) that has the same *channel* as the trace. The *size*

244     attribute of an allele is computed from its *scan*, and its *name* is independent of its *size*. The

245     *additional* attribute tells whether the allele represents an additional fragment (see main text

246     section "**Erreur ! Source du renvoi introuvable.**"). An allele does not use the *offset* attribute

247     that it inherits from the LadderFragment class.

248     A *Genotype* object regroups alleles that were found in an individual (chromatogram) at a

249     marker. It therefore points to these objects using relationships called *alleles*, *sample* and

250     *marker* respectively. A genotype stores the *a* and *b* parameters used to correct for allele

251     sizes (see main text section "**Erreur ! Source du renvoi introuvable.**") in an attribute called

252     *offset*. Indeed, these offset parameters are specific to an individual analyzed at a marker,

253     hence of a genotype. The genotype *status* attribute tells the user whether the genotype has

254     been called, whether alleles were found, the genotype has been manually edited, etc.

255     The Folder class lets users organize chromatograms and marker panels. A folder has a name

256     and points to subfolders via a set called subfolders. Reciprocally, each subfolder points to its

257     parent folder via its parent relationship. A nested hierarchy of folders can therefore be

258    constructed. As a subclass of Folder, a Panel has a name and a parent, which is an object of

259    the PanelFolder class. The subfolders of such object may be panels and/or other panel

260    folders. As a Folder, a Panel can technically have subfolders, but the application code

261    prevents this.

262    A SampleFolder is a Folder that can contains Chromatogram objects in a set called samples.

263    A SmartFolder specifies a search predicate in its *searchPredicate* attribute. When it is

264    accessed, a smart folder returns chromatogram objects found across the whole database via

265    the search predicate. A smart folder has a parent that is the same for all smart folders, it but

266    is not allowed to contain subfolders.

267

268

269

270

| Folder | A container (abstract class) |
|---|---|
| name | String |
| subFolders | Set of Folder objects |
| parent | Folder |

| SampleFolder (Folder) | A Folder containing chromatograms |
|---|---|
| samples | Set of Chromatogram objects |

| SmartFolder (Folder) | A Folder that returns chromatograms meeting search criteria |
|---|---|
| searchPredicate | Predicate |

| PanelFolder (Folder) | A Folder that contains marker panels |
|---|---|

| Panel (Folder) | A panel of STR markers |
|---|---|
| markers | Set of Marker objects |
| samples | Set of Chromatogram objects |

| Region | A named range defined in base pairs (abstract class) |
|---|---|
| name | String |
| start | Float |
| end | Float |

| Marker (Region) | A Region defining a microsatellite marker |
|---|---|
| channel | Integer (0 to 3) |
| motiveLength | Integer (2 to 7) |
| ploidy | Integer (1 or 2) |
| panel | Panel |
| bins | Set of Bin objects |
| genotypes | Set of Genotype objects |

| Bin (Region) | A Region defining a bin |
|---|---|
| marker | Marker |

| Chromatogram | Contains medatada from a chromatogram file |
|---|---|
| sampleName | String |
| well | String |
| plate | String |
| owner | String |
| runStopTime | Date |
| nScans | Integer |
| offScaleScans | Array of integers |
| coefs | Array of floats |
| … | |
| folder | SampleFolder |
| panel | Panel |
| genotypes | Set of Genotype objects |
| traces | Set of Trace objects |
| sizeStandard | SizeStandard |

| Trace | Fluorescence data from a Chromatogram at a channel |
|---|---|
| channel | Integer (0 to 4) |
| rawData | Array of integers |
| peaks | Array of structures |
| dyeName | String |
| isLadder | Boolean |
| chromatogram | Chromatogram |
| fragments | Set of LadderFragment objects |

| Genotype | Regroups Alleles found in a Chromatogram at a Marker |
|---|---|
| notes | String |
| offset | structure: *a* and *b* correction coefficients (floats) |
| status | Integer |
| sample | Chromatogram |
| alleles | Set of Allele objects |
| marker | Marker |

| SizeStandard | A standard to which a DNA ladder conforms |
|---|---|
| name | String |
| editable | Boolean |
| samples | Set of Chromatogram objects |
| sizes | Set of SizeStandardSize objects |

| SizeStandardSize | A size specified by a size standard |
|---|---|
| size | Integer |
| sizeStandard | SizeStandard |

| LadderFragment | A DNA fragment identified in a trace (molecular ladder) |
|---|---|
| name | String |
| size | Float |
| scan | Integer |
| offset | Float |
| trace | Trace |

| Allele (LadderFragment) | A DNA fragment identified in a Trace within a Marker range |
|---|---|
| additional | Boolean |
| genotype | Genotype |

271 **Figure S2.** Overview of the database managed by STRyper. Each table describes a class defined in the application. The top left cell shows the class name
272 followed by the name of its superclass in parentheses, if relevant. The top-right cell shows a brief description of the class. Below the table header, the left
273 column lists the names of attributes and relationships of the class. Relationship names are italicized. The right column specifies the type of each
274 attribute/relationship. A class inherits attributes and relationships from its superclass, but it may not use them in the application code. All relationships are
275 reciprocal, and reciprocity is represented by arrows. Single arrows point to to-one relationships (pointers to a single object) and double arrows point to to-

11

276    many relationships (sets of pointers to several objects). Only attributes and relationships saved in the database are shown. For a complete definition of these

277    classes, see header files at https://github.com/jeanlain/STRyper/tree/main/STRyper/Entities/