



Automation of high volume MPS mixture interpretation using CaseSolver

Øyvind Bleka^{a,*}, Rebecca Just^b, Jennifer Le^b, Peter Gill^{a,c}

^a Department of Forensic Sciences, Oslo University Hospital, Norway

^b DNA Support Unit, Federal Bureau of Investigation Laboratory, USA

^c Department of Clinical Medicine, University of Oslo, Norway

ARTICLE INFO

Keywords:

MPS
STR
LUS
LR

ABSTRACT

Sixty samples were sequenced using the ForenSeq kit to produce complex STR DNA mixtures, represented with three separate formats capturing different degrees of sequence information. All mixtures were run through the (open-source) CaseSolver software for comparison against both 10 reference profiles. By comparing the performance of the qualitative and the quantitative models for the different allele formats we found the following preliminary results: The quantitative model performed better than the qualitative model, however the gain was only substantial when the LR was already large. The performance gain of using the longest uninterrupted stretch over using only repeat units was large, whereas the further gain of also using the whole sequence information was small.

1. Introduction

For very serious crimes, reporting scientists often have to contend with complex investigations, where hundreds of items may be submitted by investigators for analysis per individual case. To expedite the process of comparing reference profiles to evidence profiles, many of which may be mixtures, the open-source software 'CaseSolver' was developed (www.euroformix.com/casesolver; [1]). In addition to handling size-based Short Tandem Repeat (STR) data, the software can also interpret Massive Parallel Sequencing (MPS) data to utilize sequence-level information. Here, we provide a demonstration using 60 mixtures typed using the ForenSeq DNA Signature Prep kit (Verogen, Inc., San Diego, CA). A comparison of results using different levels of sequence information was carried out using ordinary repeat units (RU), the longest uninterrupted stretch (LUS) representations [2], and LUS + representations [3] that captured all sequence variation present in the donor genotypes. We report preliminary results on the extent of enhanced discrimination offered by MPS when complex mixtures are analysed.

2. Methods

Three sets of ten samples each consisting of two, three and four-person mixtures were typed in duplicate. Contributor ratios ranged from 2:1 to 40:1 and 20:5:5:1. Donor templates ranged from 1.463-0.043 ng for the major contributor, and 0.250-0.002 ng for minor

contributors. Total DNA inputs were 1.5-0.067 ng. Libraries were pooled in sets of 32 for MiSeq FGx sequencing. Data analyses in the ForenSeq Universal Analysis Software (UAS; Verogen, Inc.) were performed using an analytical threshold (AT) of 4.5%, and the stutter thresholds taken from [4]. The lookup tables from [3] were used to translate UAS sequence strings to LUS and LUS + alleles.

CaseSolver (v1.5.0, with euroformix_2.2.0) was used to obtain likelihood ratio (LR) values (using both qualitative and quantitative models) for all 600 comparisons (by setting allele matching and LR thresholds to zero) for each allele representation format (RU, LUS and LUS +). The settings were AT = 30, default drop-in model, and stutter and degradation models turned off. The allele frequencies used were the African-American population from [3], developed using the sequence frequencies from [5].

3. Results and discussion

Application of the 4.5% AT eliminated all non-allelic sequences in the mixtures. However, with use of a dynamic AT, many authentic alleles from low-level minor donors were also removed when donor ratios differed greatly and total locus reads exceeded 650 due to the major donor. This contributed to under-estimation of the contributor number in 20.6% of 'Hp true' tests, and instances of false exclusion of some donors in these scenarios (see Fig. 1). As would be expected, inconclusive LRs also occurred with very low donor contributions.

We first compared the LRs (log₁₀ scale) of the qualitative and the

* Corresponding author.

E-mail address: oyvble@hotmail.com (Ø. Bleka).

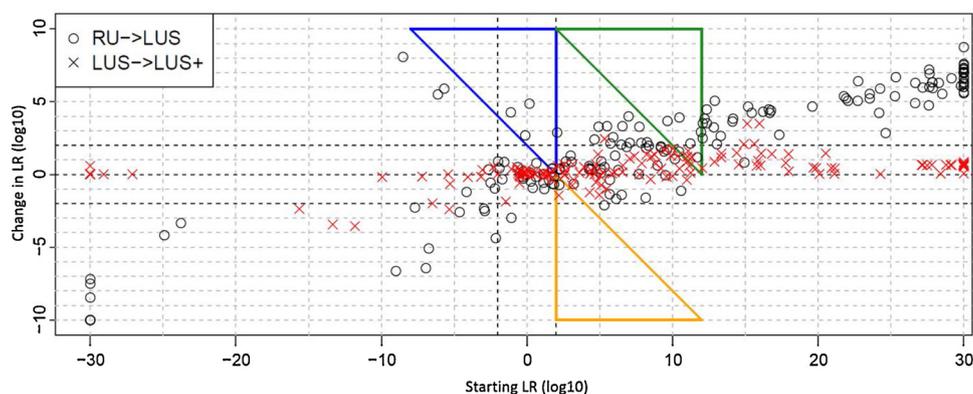


Fig. 1. LR change for Hp true tests with use of more sequence information for quantitative model. The black circles reference the LR change from RU to LUS alleles, while the red Xs indicate the LR change from LUS to LUS + alleles (log10 scale). The vertical dashed lines indicate starting LR (log10) that are inconclusive (-2 to 2), provide support for inclusion (2–12) and provide very strong support for inclusion (> 12). The green triangle indicates some support for inclusion initially (RU or LUS) but provided very strong support for inclusion when more sequence information (LUS or LUS +) was used. The blue triangle indicates an inconclusive LR (or support of Hd) initially but provided support for inclusion when more

sequence information was used. The orange triangle indicates some support for inclusion initially, but produced an inconclusive LR (or support of Hd) when more sequence information was used. Points are truncated to [-10,10] on y-axis and [-30,30] on x-axis.

quantitative models for the three different formats; RU/LUS/ LUS + . When the qualitative LR were in the range 2–12, the gains achieved by use of the quantitative model were modest: the LR increased by > 1 order of magnitude 39/32/34% of the time, and the average LR increase was 1.0/1.1/1.1. When qualitative LR exceeded 12, the increased quantitative LR were both more variable and typically much larger.

The difference in LR by use of the qualitative vs quantitative model was considered with respect to the extent of sequence information utilized for interpretation. Though not shown here, the results indicated: 1) the shift from using RU alleles to LUS alleles produced greater changes in the LR than were observed with LUS vs. LUS + interpretations, and 2) the changes in the LR for the quantitative vs qualitative models were more substantial for RU→LUS than for LUS→LUS + . Three outlier data points indicated a multiple order of magnitude change in the LR between the RU and LUS interpretations under the quantitative model, but no similar change under the qualitative model. These occurred when the LUS increased the estimated number of contributors (one more compared to the RU) such that the quantitative model better estimated the major.

Fig. 1 displays the change in LR for 180 Hp true hypotheses as progressively more sequence information was utilized for interpretation (RU→LUS, and LUS→LUS+) under the quantitative model. Overall, greater change in the LR for RU→LUS as compared to LUS→LUS + was observed. When LR (log10) was inconclusive (-2 to 2) for RU, use of the LUS changed the LR (log10) by 0.3 on average, while from LUS to LUS + the average LR change was only -0.02. When the RU LR (log10) provided some support for inclusion (2–12), the increase in LR (log10) by use of LUS averaged 1.2 (3.6 if any support), and when the LUS LR provided any support for inclusion (> 2), whereas the increase in LR (log10) by use of LUS + averaged 0.6. Similarly, when starting LR (log10) values were in the range of 2–12, the LR increased by 2+ orders of magnitude more frequently when LUS was compared to RU interpretations (31%) versus when LUS + was compared to LUS interpretations (4%). The green triangle indicates the interpretation change from some support for Hp (using RU or LUS) to very strong support for Hp when more sequence information was used (LUS or LUS+). Here the number of data points for RU-> LUS was six, versus three for LUS→LUS + . Similarly, the blue triangle indicates the change from no support for Hp to support for Hp. Here the number of points for RU-> LUS was five, versus zero for LUS→LUS + . The orange triangle indicates the change from Hp support to no support for Hd. Here the number of points for both RU-> LUS and LUS→LUS + was one.

4. Conclusion

We have demonstrated that CaseSolver can be used to perform large scale comparisons of MPS datasets, including highly complicated profiles. The preliminary data described here also clearly indicated greater LR gains when comparing LUS to RU alleles than when comparing LUS + to LUS alleles. This was due to the fact that the LUS designations captured the majority of the sequence variation present in the donor genotypes; as a result, use of the LUS + designations produced few additional alleles that could result in substantial changes to the LR. The data produced in this study also highlighted what impact the use of a dynamic AT may have on allele recovery, which in turn influences downstream probabilistic interpretation. Further work will be carried out to reanalyse mixtures using a static AT, with subsequent re-interpretation with CaseSolver.

Declaration of Competing Interest

The authors declare no conflicts of interest.

Acknowledgements

The authors thank Jodi Irwin and Anthony Onorato of the FBI Laboratory. This research was supported in part through the FBI's Visiting Scientist Program, an educational opportunity administered by the Oak Ridge Institute for Science and Education (ORISE), and by an ISFG Short-Term Travel Fellowship. Names of commercial manufacturers are provided for identification purposes only, and inclusion does not imply endorsement of the manufacturer, or its products or services by the FBI. The views expressed are those of the authors and do not necessarily reflect the official policy or position of the FBI or the U.S. Government.

References

- [1] Ø Bleka, L. Prieto, P. Gill, CaseSolver: an investigative open source expert system based on EuroForMix, *Forensic Sci. Int. Genet.* 41 (2019) 83–92.
- [2] R.S. Just, J.A. Irwin, Use of the LUS in sequence allele designations to facilitate probabilistic genotyping of NGS-based STR typing results, *Forensic Sci. Int. Genet.* 34 (2018) 197–205.
- [3] R.S. Just, J. Le, J.A. Irwin. LUS+: Extension of the LUS designator concept to differentiate most sequence alleles for 27 STR loci Submitted to *Forensic. Sci. Int. Genet.*
- [4] L.I. Moreno, M.B. Galusha, R. Just, A closer look at Verogen's Forenseq DNA Signature Prep kit autosomal and Y-STR data for streamlined analysis of routine reference samples, *Electrophoresis* 39 (2018) 2685–2693.
- [5] K.B. Gettings, et al., Sequence-based U.S. population data for 27 autosomal STR loci, *Forensic Sci. Int. Genet.* 37 (2018) 106–115.