

1 **Title:**

2 Alternative LC-MS/MS Platforms and Data Acquisition Strategies for Proteomic Genotyping of Human
3 Hair Shafts

4
5 **Authors:**

6 Zachary C. Goecker¹, Kevin M. Legg², Michelle R. Salemi³, Anthony W. Herren³, Brett S. Phinney³,
7 Heather E. McKiernan², Glendon J. Parker^{1,*}

8 ¹ Department of Environmental Toxicology, University of California, Davis, CA, USA

9 ² The Center for Forensic Science Research and Education, Willow Grove, PA, USA

10 ³ Proteomics Core Facility, University of California, Davis, CA, USA

11 * Corresponding author.

12 Glendon Parker PhD

13 Department of Environmental Toxicology

14 University of California – Davis

15 One Shields Ave. Davis, California 95616

16 p. (530) 752 9870

17 e. gjparker@ucdavis.edu

18

19 **Running Title:** Alternative Platforms for Proteomic Genotyping

20

21

22

23

24

Alternative Platforms for Proteomic Genotyping

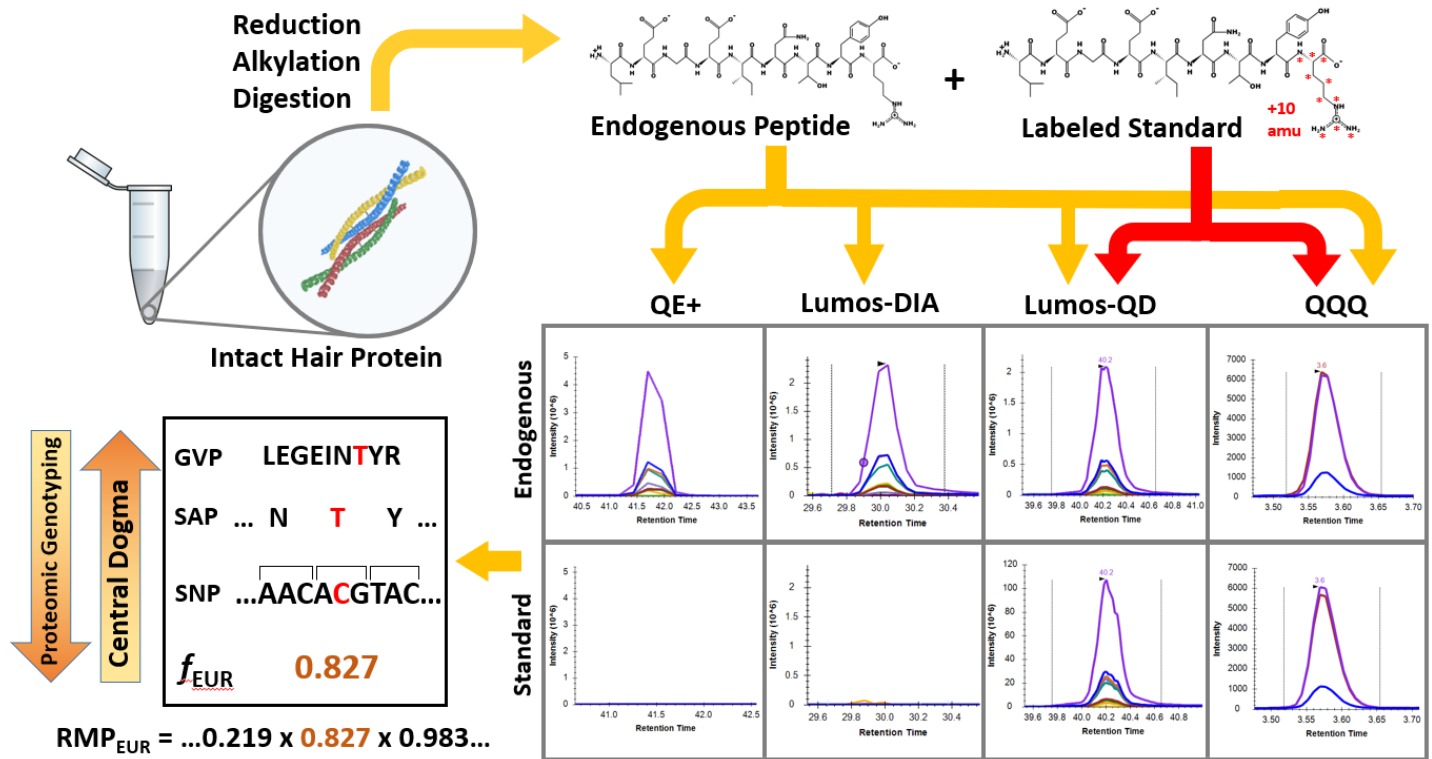
25 Highlights:

- 26
- Test four mass spectrometry configurations to optimize detection of genetically
- 27 variant peptides
- Technology transfer of proteomic genotyping assays
- 28
- Improved sensitivity results in higher level of forensic discrimination for human
- 29 identification using multiple reaction monitoring
- 30

31

32 Graphical Abstract:

33



35

36

37

Alternative Platforms for Proteomic Genotyping

38 **Abstract:**

39 Protein is a major component of all biological evidence. Proteomic genotyping is the use of genetically
40 variant peptides that contain single amino acid polymorphisms to infer the genotype of matching non-
41 synonymous single nucleotide polymorphisms for the individual who originated the protein sample. This
42 can be used to statistically associate an individual to evidence found at a crime scene. The utility of the
43 inferred genotype increases as the detection of genetically variant peptides increases, which is the direct
44 result of technology transfer to mass spectrometry platforms typically available. Digests of single (2 cm)
45 human hair shafts from three European and two African subjects were analyzed using data dependent
46 acquisition on a Q-Exactive™ Plus Hybrid Quadrupole-Orbitrap™ system, data independent acquisition
47 and a variant of parallel reaction monitoring on a Orbitrap Fusion™ Lumos™ Tribrid™ system, and
48 multiple reaction monitoring on an Agilent 6495 triple quadrupole system. In our hands, average
49 genetically variant peptide detection from a selected 24 genetically variant peptide panel increased
50 from 6.5 ± 1.1 and 3.1 ± 0.8 using data dependent and independent acquisition to 9.5 ± 0.7 and $11.7 \pm$
51 1.7 using parallel reaction monitoring and multiple reaction monitoring ($p < 0.05$). Parallel reaction
52 monitoring resulted in a 1.3-fold increase in detection sensitivity, and multiple reaction monitoring
53 resulted in a 1.6-fold increase in detection sensitivity. This increase in biomarker detection has a
54 functional impact on the statistical association of a protein sample and an individual. Increased
55 biomarker sensitivity, using Markov Chain Monte Carlo modeling, produced a median estimated random
56 match probability of over 1 in 10 trillion from a single hair using targeted proteomics. For parallel
57 reaction monitoring and multiple reaction monitoring, detected genetically variant peptides were
58 validated by the inclusion of stable isotope labeled peptides in each sample, which served also as a
59 detection trigger. This research accomplishes two aims: the demonstration of utility for alternative
60 analytical platforms in proteomic genotyping, and the establishment of validation methods for the
61 evaluation of inferred genotypes.

Alternative Platforms for Proteomic Genotyping

62 **Abbreviations:** DDA, data dependent acquisition; DIA, data independent acquisition; GVP, genetically
63 variant peptide; HID, human identification; MRM, multiple reaction monitoring; PRM, parallel reaction
64 monitoring; QD, QuanDirect; QQQ, triple quadrupole; RMP, random match probability; SAP, single
65 amino acid polymorphism; SIL, stable isotope labeled; SNP, single nucleotide polymorphism

66

67 **Keywords:** Hair, Forensic Proteomics, Genetically Variant Peptides, Human Identification, Proteomic
68 Genotyping

69

70 **1. Introduction**

71 Proteomics has many promising applications in a legal context, with recent advances being
72 made in body fluid identification, drug interactions, sex estimation, and human identification (HID) [1–
73 6]. Proteomic genotyping is the analysis of protein sequence variants, termed single amino acid
74 polymorphisms (SAPs), to infer single nucleotide polymorphism (SNP) alleles. SNPs can be inherited and
75 subsequently detected in peptides from digested protein [3]. These peptides, termed genetically variant
76 peptides (GVPs), are especially useful in samples where DNA may be absent or highly degraded, such as
77 is often the case with hair, fingerprints, and bone [7, 8]. Currently, mitochondrial DNA is the mainstay in
78 HID for highly degraded samples and archaeological remains due to its multiplicity in each cell, rapid
79 evolution, and familial information [9, 10]. Proteomic genotyping offers intrinsic advantages compared
80 to other DNA-based genetic analyses. Protein, much like mitochondrial DNA, has many copies per cell.
81 Processing of protein does not involve an amplification step or preliminary knowledge of the sequence,
82 as is the case for DNA primer design. The peptide bond is chemically stable and common chemical
83 modifications are predictable and accommodated by spectral matching algorithms [11, 12]. Protein has
84 proven to outlast DNA due to degradation effects for archaeological sex estimation from teeth [6].
85 Tryptic peptides average only 14 residues, which reduces the probability of random cleavage and

Alternative Platforms for Proteomic Genotyping

86 information loss [13]. The stoichiometry of protein copy number is greater than DNA by up to 7 orders
87 of magnitude with a median increase of 5 orders of magnitude, allowing for detection and analysis
88 without amplification [14–16].

89 Initial development of proteomic genotyping has focused on hair shafts as a source of protein-
90 based genetic information. Over 400 genetically variant peptides from hair have been identified that can
91 predict the corresponding SNP allele [17]. Optimization to this point has focused on chemical processing
92 to maximize peptide production from a single hair shaft [17–19]. Aside from a shift to more modern
93 mass spectrometry instruments for data dependent acquisition, relatively little has been done to
94 optimize data acquisition at the mass spectrometry level. Another side effect of the current focus on
95 data dependent acquisition for proteomic genotyping is the dependence on peptide spectral matching
96 for detection of genetically variant peptides. Peptide spectral matches, depending on the algorithm,
97 come with a statistical expectation score or other measures of confidence [11, 12, 20–22]. For most
98 proteomic applications this is not functional, since multiple peptides are often identified for each gene
99 product and the level of uncertainty can be miniscule [22]. Proteomic genotyping, however, relies on
100 single peptides to infer SNP alleles. Validation of the inferred SNPs are therefore necessary and most
101 easily provided by direct confirmation of genotype by DNA sequencing. Validation is also possible to
102 confirm peptide identification through the addition of stable isotope labelled (SIL) synthetic peptides
103 into a sample digest. These standards, equivalent in sequence and chemistry to the matching
104 endogenous peptides, behave identically to matching endogenous peptides, and do not interfere with
105 endogenous peptide detection. The use of SIL peptides is a standard feature of targeted mass
106 spectrometry platforms for use in triggering endogenous peptide detection and quantification [23].

107 For proteomic genotyping to be readily available to forensic investigators, it also needs to be
108 conducted on platforms that are widely accessible to investigators. Targeted mass spectrometry using
109 triple quadrupole systems is commonly used for many forensic toxicology analyses. It is also affordable,

Alternative Platforms for Proteomic Genotyping

110 robust, and reproducible for both chromatography and mass spectrometry. Use of targeted methods of
111 mass spectrometry potentially improves sensitivity and bolsters analyte identification confidence,
112 helping to fulfill the guidelines in forensic analyte identification. Guidelines set by the scientific working
113 group for forensic toxicology (SWGTOX) [24], European Commission (directive 96/23/EC) [25], and
114 World Anti-Doping Agency (WADA) [26] all list a minimum number of identification points to confirm the
115 presence of a drug or other analyte. These identification points are derived from retention time
116 windows, peak shape, and transitions and may not be satisfied using standard proteomic discovery
117 techniques such as data dependent acquisition (DDA) alone without DNA-based verification. Currently,
118 proteomic genotyping in forensic science has focused on the optimization of peptide production in
119 sample preparation, and expansion to other forensically relevant tissue sources [17, 19]. These studies
120 have relied on shotgun mass spectrometry and nano liquid chromatography coupled with orbitrap mass
121 spectrometry and are yet to exploit useful alternative instrumental strategies available [17, 27, 28].
122 Here, we propose spiking SIL GVPs into hair protein digests as a means of peptide identification
123 validation and as a mass trigger for data acquisition. This technique has typically been used for the
124 quantification of other peptide targets [29]. A standard will elute chromatographically and be analyzed
125 in the same time window as its corresponding endogenous GVP. Therefore, the standard and expected
126 endogenous peptides can be directly compared in terms of retention time and ion ratios (Figures S1 and
127 S2) providing a means for real time validation.

128 This study explores the capability and utility of alternative mass spectrometry platforms and
129 data acquisition strategies using a panel of GVPs. Instead of a direct comparison of the platforms and
130 acquisition methods which would involve an in-depth evaluation of instrument components, a proof of
131 concept and evaluation of performance was studied for proteomic genotyping purposes only. Two
132 approaches are first assessed: shotgun proteomics (DDA) on a Q Exactive™ Plus Hybrid Quadrupole-
133 Orbitrap™ Mass Spectrometer, and data independent acquisition (DIA) on a Orbitrap Fusion™ Lumos™

Alternative Platforms for Proteomic Genotyping

134 Tribrid™ Mass Spectrometer [30, 31]. Two other approaches were also tested in the presence of
135 matching synthetic SIL peptide standards. These include a variant of parallel reaction monitoring (PRM),
136 called QuanDirect™ (QD) [32] conducted on the tribrid system that uses detection of SIL standard
137 peptides to trigger data acquisition in the mass window of the corresponding endogenous peptide, and
138 multiple-reaction monitoring (MRM) conducted on an coupled Agilent 1290/6495 triple quadrupole
139 system [33–37]. To provide a direct comparison in the performance of genetically variant peptide (GVP)
140 detection between the four acquisition methods, hair from European-American (n = 3) and African-
141 American (n = 2) subjects was tested in three biological replicates. Results reported here are limited by
142 the procedures as performed using standard mass spectrometry platforms and protocols within the
143 range of what is considered best practice for each method. The noted enhancements, based on a panel
144 of 24 standard peptides, have the potential to dramatically increase both the discriminatory power of
145 proteomic genotyping and the applicability of the method since it uses instruments and analyte
146 detection criteria commonly found in toxicology laboratories.

147

148 **2. Experimental Section**

149 *2.1 General Experimental Design*

150 Hair shafts were collected from 5 individuals who are representative of two populous ancestral
151 backgrounds in of the United States: European and African. The number of individuals needed for this
152 study was minimal since this was a novel proof of concept study to demonstrate the usage of targeted
153 proteomics for proteomic genotyping. Enough donors were used to assess reproducibility and calculate
154 standard deviation. Three single hairs from each individual, 2 cm in length, were processed separately
155 using a previously developed method, with a total of 15 hair digests. A blank with trypsin and without
156 trypsin were also processed in parallel with all other digests. 24 stable isotope labeled (SIL) standard
157 genetically variant peptides (GVPs) were spiked into the hair digests for parallel reaction monitoring and

Alternative Platforms for Proteomic Genotyping

158 multiple reaction monitoring (Table S1). Raw mass spectral data were processed using the Skyline
159 software for the 24 GVPs of interest and their corresponding heavy-isotope peptide standards. 24
160 peptides were chosen to adequately represent the diversity of the full set of 408 currently identified
161 GVPs in terms of detection sensitivity, length, and composition. Basic statistical analyses were
162 conducted such as standard deviation, random match probability, false discovery rate ($FP/(FP+TP)$), and
163 detection sensitivity ($TP/(TP+FN)$) to compare the three analytical methodologies. Random match
164 probability calculations were estimated using the procedure outlined in Parker et al [3].

165

166 *2.2 Hair Collection and Processing*

167 Samples used in this study were prepared as part of an earlier study [17]. Briefly, five individuals
168 were analyzed: three subjects of European (Davis, CA) and two subjects of African genetic background,
169 respectively (Sorenson Forensics LLC, Salt Lake City, UT). Hair and saliva were collected using protocols
170 compliant with the Institutional Review Board at the University of California, Davis (IRB# 832726). Hairs
171 were collected by cutting a few inches inward from the distal end, therefore excluding the roots. The
172 length of hair on the head before cutting was roughly 10 cm. Hair shafts were further cut to a length of
173 20 mm before continuing with protein extraction [17]. The African hair samples weighed almost half of
174 the weight of the European hair samples due to differences in hair shaft width and shape (data not
175 shown). Hair shafts were biochemically processed using an optimized processing protocol as part of an
176 earlier study [17]. After initial preparation and use of the samples to generate the data-dependent
177 acquisition datasets used in the cited study, the remaining supernatants were stored at -20°C. Prior to
178 mass spectrometric analysis, the samples were again centrifuged to minimize insoluble particulates in
179 the supernatant.

180

181 *2.3 Selection of a Panel of Stable Isotope Labeled Genetically Variant Peptides*

Alternative Platforms for Proteomic Genotyping

182 A panel of 24 highly characterized GVPs from 12 loci was selected to represent a wide range of
183 potential sensitivities, from rarely to frequently detected, and were well characterized over a range of
184 studies and laboratory groups [3, 17, 19, 27, 38]. The peptides selected range from 8 to 21 amino acids
185 in length and were all modified with stable isotopes at the C-terminal lysine (+8 Da) or arginine (+10 Da)
186 (JPT peptide technologies, Acton, MA) (Table 1). The peptides (14 nmol/well) were subsequently
187 suspended in 4 μ L of 70% formic acid and 136 μ L of 0.1% trifluoroacetic acid. 10 μ L of each standard
188 were pooled to make a concentration of 4.16 pmol/ μ L per standard. The pooled sample was then
189 purified using a silica C18 macrospin column with loading capacity of 30-300 μ g of peptide material (The
190 Nest Group, Southborough, MA). Briefly, the peptide digests were loaded onto the column and spun,
191 the column was washed with 0.1% trifluoroacetic acid and spun three times, and the peptides were
192 eluted using 80% acetonitrile 0.5% formic acid and spinning and were subsequently dried down. This
193 pool was then injected into four hair digests from two individuals at 1 and 2 fmol/ μ L. A second pool was
194 then created, with normalization based on peak area to make a final spike mixture (Table S1). This final
195 spike solution consisted of a total of 433 nmol of peptides in 1 mL final volume (433 pmol/ μ L). The final
196 mixture (1 fmol) was included in each sample injection for QuanDirect analysis and 73 pmol of the final
197 mixture was injected into each sample for MRM analysis on the triple quadrupole (QQQ) platform.

198

199 *2.4 Instrumental Analysis*

200 Before applying the samples to LC-MS/MS, solubilized tryptic peptides were quantified using the
201 Pierce™ Quantitative Fluorometric Peptide Assay (ThermoFisher) as reported in previous work [17].
202 Resulting data were used to determine how much material to apply to the instrument. For the Q
203 Exactive Plus and Fusion Lumos, this amounted to 0.75 μ g of peptide digest material. Digest injection
204 volume was held constant on the QQQ platform, but volumes varied based on concentration for the
205 other three platforms.

Alternative Platforms for Proteomic Genotyping

206 Three instruments were used to conduct four data acquisition methods (Figure 1). The first
207 analysis method (QE+) was conducted on a Q Exactive Plus nLC-MS/MS platform which employed data
208 dependent acquisition (DDA). This method was established as part of an earlier study [17]. For the
209 second analysis (Lumos-DIA), samples were analyzed on a Thermo Scientific Fusion Lumos mass
210 spectrometer and was connected to a Dionex nano Ultimate 3000 (Thermo Scientific) with a Thermo
211 Easy-Spray source. The acquisition method was set to data independent acquisition (DIA). For this
212 method, peptides were trapped and separated on a 100 μm x 250 mm C18 column with 3 μm particle
213 size PepMap Easy-Spray (Thermo Scientific) using a Dionex Ultimate 3000 nUPLC at 200nl/min. Peptides
214 were eluted using a 90 min gradient of 0.1% formic acid (A) and 80% acetonitrile (B). Gradient conditions
215 include 2% B to 50% B over 60 minutes, followed by a 50%-99% B in 6 minutes and then held for 3
216 minutes, then 99% B to 2% B in 2 minutes. The mass spectrometer was run in DIA mode using a collision
217 energy of 35, resolution of 30K, maximum inject time of 54 ms and an automatic gain control (AGC)
218 target of 50,000. Each individual sample was run in DIA mode with staggered isolation
219 windows of 12 Da in the range 400-1000 m/z. For each analytic sample, the individual sample was run in
220 DIA mode using the same settings as the chromatogram library runs except using staggered isolation
221 windows of 8 Da in the m/z range 400-1000 m/z.

222 The third analysis method (Lumos-QD) was conducted on the same instrument as the Lumos-
223 DIA method, except QuanDirect™ (QD) parallel reaction monitoring acquisition was employed. For this
224 method, the digested peptides were reconstituted in 2% acetonitrile/0.1% trifluoroacetic acid and 1 μg
225 in 5 μl of each sample was loaded onto a 75 μm x 20 mm PepMap 100Å 3U trap (Thermo Fisher
226 Scientific) where they were desalted online before being separated on a 50 μm x 150 mm 100 Å 2U
227 PepMap EasySpray column (Thermo Fisher Scientific). Peptides were eluted using a 90 min gradient of
228 0.1% formic acid (A) and 80% acetonitrile (B) with a flow rate of 200nL/min. Gradient conditions include
229 2% to 50% B over 66 min, 50% to 99% B over 2 min, and then held at 99% B for 4 min, followed by 99%

Alternative Platforms for Proteomic Genotyping

230 to 2% B over 2 min. Targeted precursors were interrogated for a maximum 3 sec cycle. The target list
231 consisted of the m/z ratios and charge states of the heavy peptides (Table 1) only since no retention
232 time info was required. The QuanDirect™ method of targeted mass spectrometry was used to search for
233 endogenous GVPs [32, 39]. To trigger a data-dependent scan, the precursor must match the expected
234 charge state and the m/z within 10 ppm. These precursors were interrogated with a short SRM scan
235 (MS^2 IT HCD) of the predicted y_1 fragment ion region for the expected heavy R (180-190 m/z) or heavy K
236 (150-160 m/z) peptides. The y_1 fragment in proteomics determines the C-terminus ending and is more
237 easily identifiable in a mass spectrum due to its monomeric independence. The ultra-fast SRM scans
238 were performed using a 10 m/z mass range, rapid ion trap scan rate, HCD NCE 40%, 0.7 m/z isolation
239 window, AGC target 1E4, and a max IT of 10 ms. The y_1 ion in the SRM scan must be above an intensity
240 threshold of 1000 and within 1 Da of the expected m/z . If the expected heavy y_1 fragment ion was
241 detected in the SRM scan, 185.1 (heavy R) or 155.1 (heavy K), then full HRAM HCD MS/MS scans were
242 triggered on the spiked-in heavy peptide as well as the endogenous form (Table S2). These scans use the
243 following parameters: scan range 150-1500 m/z , 60K resolution, 30% NCE HCD, AGC target 2e5, and max
244 IT of 110 ms. To trigger on the endogenous peptide, an m/z offset of -5 or -4 U was used for the R and K
245 peptides, respectively.

246 For the fourth analysis (QQQ), samples were analyzed via multiple reaction monitoring
247 acquisition using an Agilent 1290 Infinity series HPLC system, coupled to an Agilent 6495 triple
248 quadrupole mass spectrometer with an Agilent Jet Stream source. 10 μ L of a 1:10 spike:digest (v/v) ratio
249 (~15 μ g digested material and 291 pg of spike) was loaded on a 2.1 mm \times 100 mm, 2.7 μ m AdvanceBio
250 Peptide Map fused-core silica column (Agilent), and separated over a 15 min gradient at 400 μ L/min.
251 The solvent gradient for the elution of peptides began with 5% ACN and increased to 35% ACN at 11
252 min, 65% ACN at 12.5 min, and 90% ACN at 13 min and held for 2 min, and then reduced to 5% for 5 min
253 to re-equilibrate the column. Source conditions included a gas temperature at 150°C at a flow rate of 11

Alternative Platforms for Proteomic Genotyping

254 l/min, nebulizer pressure of 30 psi, sheath gas temperature of 150°C at 10 l/min, and a capillary voltage
255 of 3500 V. Collision energies were calculated based on precursor m/z and charge state in Skyline
256 software, and were not fully optimized. Data were acquired in positive dynamic MRM mode (dMRM)
257 with an MS^1 resolution set to wide and MS^2 resolution to unit, retention time window of 30 sec, and a
258 cycle time of 500 ms. Three transitions were selected for the detection of each standard and
259 endogenous peptide (Table S3).

260

261 *2.5 Software Analysis*

262 To make the resulting QuanDirect™ PRM datafiles amenable with Skyline software, spectra from
263 the linear ion trap were excised using the FT RecalOffline tool from Xcalibur™ (ThermoFisher Inc.). This
264 treatment does not interfere with spectrum interpretation since this was only a part of the internal
265 decision tree. Raw files were manually loaded and the external slicer was called, under Rawfile
266 Functions, from within RecalOffline to remove any masses below 200 m/z using a mass filter from 200
267 m/z to 2000 m/z . Since only low mass y_1 fragments of 185 and 155 m/z with scan range < 200 m/z were
268 searched for, the filter removed all of the ion trap data from the file. The resulting file only contained
269 MS^1 and MS^2 orbitrap data which was used to analyze endogenous and standard peptides.

270 Skyline software [40] (version 20.1) was used to visualize endogenous peptide data and SIL
271 peptide data simultaneously and to extract only mass transitions of interest (Tables S2 and S3). Positive
272 peptide identification required a s/n ratio > 3, peak intensity of ion targets > 20 counts, and an ion ratio
273 between quantifier and qualifier ions within 25% of the target. These parameters were chosen to meet
274 minimum identification points from common forensic guidelines [24–26]. For the analysis performed
275 here, samples were analyzed in batches based on the instrumentation on which they were run. For
276 samples analyzed on the QE+ using DDA, full-scan transition settings for MS^1 filtering were set to include
277 count isotope peaks, orbitrap precursor mass analyzer, with 3 peaks and a resolving power set to 60,000

Alternative Platforms for Proteomic Genotyping

278 at 400 m/z . MS/MS filtering settings were set to DDA as the acquisition method, orbitrap product mass
279 analyzer, no isolation scheme, a resolving power set to 60,000 at 400 m/z . All other settings were set to
280 default. For samples analyzed using Lumos-DIA, the same settings were used as DDA with the exception
281 of 70,000 MS^1 resolving power, DIA as the acquisition method, results only as the isolation scheme, and
282 17,500 as the MS/MS resolving power. For samples analyzed using Lumos-QD, the same settings were
283 used as Lumos-DIA except for no isotope peaks, MS/MS filtering settings were set to targeted as the
284 acquisition method, and no isolation scheme. For samples analyzed using QQQ with MRM, the same
285 settings were used as PRM except for MS^1 filtering were set to include count isotope peaks. One
286 precursor and 3 transitions were chosen for MRM analysis, while one precursor and 10-24 transitions
287 were chosen for the PRM analysis and DDA analysis. These are both above the minimum standard
288 guideline for the number of ions required for a positive identification [25, 26].

289 Positive peptide identifications were called from the Skyline software based on precursor and
290 transition signal to noise ratio, retention time, transition masses, and ion ratios. For retention time, this
291 identification criteria included having a GVP retention time within 2% or ± 0.1 min of the labeled
292 standard. For the DIA and DDA approaches, comparison of retention time to a labeled standard was not
293 used. In terms of signal to noise ratio, a minimum ratio of 3:1 was used as the threshold for data from all
294 platforms. Transitions used for identification for the targeted approaches were taken from the most
295 abundant transitions for the labeled standard peptides. The untargeted methods utilized the Prosit
296 library [41] to compare both transitions for identification and ion ratios. For all acquisition
297 methodologies, ion ratio maximum tolerance windows were set to be within 10% of the relative
298 abundance of the compared ion, so long as the peak is at least 50% of the base peak [26]. The PRM and
299 MRM methods used the labeled standard peptides as a reference and the DDA and DIA methods used
300 the Prosit library as the reference for ion ratios.

301

Alternative Platforms for Proteomic Genotyping

302 2.6 Statistical Analysis

303 GVP Finder (v1.2) (<https://www.parkerlab.ucdavis.edu>) was used to estimate random match
304 probabilities (RMPs). This is an excel spreadsheet compatible with X!Tandem output developed in
305 previous work [17]. In short, RMP was calculated using the product rule [3, 42] by simply multiplying
306 independent genotypic frequencies based on observations on individual genotypes from the major
307 populations in the 1000 Genome Project Consortium [43]. To account for linkage disequilibrium, it was
308 assumed that there was complete linkage for GVPs shared within an open reading frame and complete
309 independence between each open reading frame [3]. For GVPs that were determined to be genetically
310 linked within an open reading frame, a cumulative genotypic frequency was calculated by counting the
311 number of individuals in the consortium who have the same gene specific GVP profile as was obtained
312 from the sample and dividing by the total number of individuals in the population. Genetic validation
313 was performed to assign trueness of positive and negative detections. Genomic DNA was extracted and
314 sequenced as reported in previous work [17].

315 To estimate random match probability of a profile that would result from a targeted QQQ
316 analysis using all known GVPs, and assuming equivalent detection sensitivity obtained from the panel, a
317 Markov Chain Monte Carlo (MCMC) model was developed. MCMC is an algorithm that simulates
318 stochastic processes such as sampling from a probability distribution [44, 45]. This method of sampling
319 allows an estimation of true population probability distributions by randomly sampling from
320 probabilistic data. For this study, MCMC was developed as a function of GVP number and validated by
321 superimposing actual RMP values from previous studies. The probability distributions were taken from
322 actual genotype frequencies from the 408 GVPs that have been identified. A theoretical genotype of
323 non-synonymous SNP alleles was generated based on randomly selecting known GVPs and randomly
324 determining if a theoretical genotype would include that GVP based off its genotype frequency. One
325 hundred iterations were included in this model. Minimum, maximum, and median theoretical RMP

Alternative Platforms for Proteomic Genotyping

326 values were estimated based on which theoretical GVPs were randomly chosen in the model. The model
327 assumes one GVP locus per open reading frame. The resulting modelled genotypes were randomly
328 selected in each iteration as a function of prior probability based on the genotype frequency chosen
329 randomly here, rather than favoring GVPs of historically higher detection. Therefore, this model is not
330 biased towards specific GVPs and does not mimic biological GVP profile distributions.

331

332 *2.7 Data Reporting and Availability*

333 All RAW data files containing detected endogenous peptides and SIL GVPs from hair digests
334 mentioned in this work, including from the supplemental section, are publicly available on
335 ProteomeXchange (PXD024651) [46]. A complete list of datafiles is also available (Table S4). Files from
336 QE+ are comprehensive and include all detected ions, whereas the Lumos-QD and QQQ files are limited
337 to ions corresponding to GVPs from the panel. Lumos-QD data are modified to exclude linear ion trap
338 data. Skyline files for data obtained from the four analytical platforms are publicly available at
339 <https://panoramaweb.org/TargetedGVP.url> [47].

340

341 **3. Results and Discussion**

342 Studies optimizing the detection of genetically variant peptides (GVPs) have done so by focusing
343 on the chemical release of tryptic peptides from the hair matrix, or by applying the resulting peptide
344 mixtures to more sensitive instrumentation. In this study different mass spectrometry data acquisition
345 methods were tested to evaluate additional options for increased GVP detection and therefore further
346 increase the utility of proteomic genotyping in forensic investigation. Accordingly, mass spectrometry
347 data acquisition using Data Independent Acquisition (DIA), Parallel Reaction Monitoring (PRM) and
348 Multiple Reaction Monitoring (MRM) were all tested on instruments with configurations that are
349 standard for each method. Acquired data from all three platforms, as well as existing data using a

Alternative Platforms for Proteomic Genotyping

350 shotgun proteomics Data Dependent Acquisition (DDA) approach, were screened for detection of a
351 panel of 24 endogenous GVPs in replicate trypsin digests using a common bioinformatic workflow in
352 Skyline (Figures 2, 3, and S1). The cumulative inferred non-synonymous SNP genotypes were directly
353 validated using the exome of each subject (Figure 3) to determine basic metrics such as true positive,
354 false positive, true negative, and false negative rates, along with sensitivity ($TP/(TP+FN)$) and false
355 discovery rates ($FP/(FP+TP)$). Besides this binary classification process, other metrics such as signal to
356 noise, peak shape, ion ratio, peptide ionization efficiency, retention time, and abundant transitions were
357 also measured. In the case of the PRM and MRM acquisition methods, the evaluation was facilitated by
358 using a panel of exogenous stable isotope labeled peptides. While each data acquisition approach was
359 within the range of normal best practice, no systematic optimization occurred beyond establishing basic
360 acquisition and chromatographic parameters. The results therefore reflect different chromatographic,
361 ionization, and mass spectrometer systems and configurations for each acquisition method. While direct
362 comparisons could not be made, the performance of each method could be individually evaluated in
363 comparison to previously acquired data using shotgun proteomics (DDA).

364

365 *3.1 Analysis of Performance for the QE+ Platform*

366 Data previously acquired on a nano-LC / Q Exactive Plus (QE+) platform with data dependent
367 acquisition (DDA) was reanalyzed using Skyline software (version 20.1) to provide a benchmark for other
368 data acquisition strategies and instrument configurations [17, 40]. The percentage of true positive
369 shotgun proteomic identifications for the QE+ platform was 26.4% and the detection sensitivity
370 ($TP/(TP+FN)$) was 43.2% (Figure 3). The false discovery rate ($FP/(FP+TP)$) was the lowest for the QE+
371 platform at 2.6% and the average GVP detection from the 24-GVP panel was 6.5 ± 1.1 (Figure 3). Ten of
372 the 24 endogenous peptides were not detected at all using this platform. These peptides include
373 GILVDTSR, ALETVQER, ALETQER, EWSTFAVGPQHCLQLNDR, GVALSNVIHK, and GVALSNVVHK from

Alternative Platforms for Proteomic Genotyping

374 proteins HEXB, GSDMA, NEU2, and SERPINB5 (Figures 3 and S3). These peptides may not have been
375 observed due to low protein abundance in the hair sample digests, whereas keratin proteins are very
376 abundant [48]. No peptides were observed using this method that were not also observed in the other
377 methods. Four endogenous peptides were not detected at all in this study regardless of the platform
378 used: VSAMYSSSSCKLPSLSPVAR, VSAMYSSSPCKLPSLSPVAR, EHCSACGPLSQLLVK, and
379 EWSTFAVGPQHCLQLHDR. These longer peptides may be more challenging to detect based on their
380 length and residue composition [49]. For a positive detection, the average signal to noise ratio was
381 higher than 1:3 (Figure 2, S5), defined as the variance of amplitude of the baseline and signal was the
382 amplitude of the peak as measured from the apex to the average baseline, which meets the minimum
383 requirements in toxicology scientific working group guidelines [24–26]. The nanoflow chromatography /
384 Q-Exactive configuration used in this analysis was sensitive [50], requiring only 1 µg of the roughly 100
385 µg protein present in 2 cm of a single hair shaft [51]. This acquisition method can be used for GVP
386 discovery and provide a resource for retrospective analysis of GVPs, although the analysis was limited to
387 the 24 GVPs and associated ions in the SIL-peptide panel. In terms of the peak shape for the QE+
388 analytical platform, the overall form differs from the three other platforms due to differences in
389 chromatography (Figure 2). Complex chromatography patterns due to internal prolines were detected
390 [52, 53].

391

392 *3.2 Analysis of Performance for the Lumos-DIA Platform*

393 For Data Independent Acquisition (DIA), samples were applied to a Dionex nano Ultimate 3000
394 coupled to an Orbitrap Fusion™ Lumos™ Tribrid™ Mass Spectrometer and data acquired using SWATH-
395 MS (DIA) (Figure 1). Overall, Lumos-DIA peaks appear sharp with peak intensities averaging at 1×10^3
396 and the average signal to noise ratio was also above 1:10 (Figure 2). The resulting percentage of true
397 positive identifications for Lumos-DIA was 13.2% and the false discovery rate was 5.0% (Figure 3).

Alternative Platforms for Proteomic Genotyping

398 Average GVP detection from the 24-GVP panel was rather low, at 3.1 ± 0.8 and therefore detection
399 sensitivity (TP/(TP+FN)) was also low, at 21.6% (Figure 3 and S3) which is less than half the sensitivity we
400 typically achieve using standard DDA methodologies on the QE+. Overall, this method did not detect 15
401 of the 20 peptides that were detected using the other methodologies. Seven of these 15 missing
402 peptides were found in all three other methodologies, which include peptides DSQECILmETEAR,
403 LEGEINmRY, EHCSACGPLSR, DLNMDCmVAEIK, DLNMDCIVAEIK, GAFLYPCGVSTPVLSTGVLR, and
404 GAFLYdPCGVSTPVLSTGVLR from KRT39, KRT32, KRT39, KRT83, and KRT82 (Figures 3 and S3). Only one
405 peptide was observed using the Lumos-DIA methods that was not observed in either PRM or MRM
406 method for the same donor; AKPLEQAVAAIVCTFQEYAGR.

407 DIA requires little method optimization [30, 54], and may be used for GVP scouting or
408 retrospective use. Unbiased detection of peptides, uniquely for DIA, allows for proteins of low
409 abundance to be detected [31]. This precludes the need for an exclusion list or other mass filtering
410 parameter optimizations. The data is highly reproducible [31] and so running evidence samples
411 alongside exemplars would be more consistent and would result in less variance due to protein
412 abundance levels [55]. The main challenge in DIA interpretation, at least in our hands, was
413 deconvolution of MS^2 spectra. In this data false positive identification occurred in four peptides that
414 were not explained by instrument carry-over or genetics. Due to the nature of SWATH mass
415 spectrometry, an MS^2 mass spectrum may contain product ions from multiple precursor ions, which may
416 lead to convoluted MS^2 extracted ion chromatograms. This may be problematic in a courtroom setting,
417 although the use of the internal standard SIL peptides would have significantly aided MS^2
418 interpretation. The sensitivity of this analysis may also be improved with better precursor validation,
419 library match validation, and staggered SWATH windows [56].

420

421 *3.3 Analysis of Performance for the Lumos-QD Platform*

Alternative Platforms for Proteomic Genotyping

422 Targeted parallel reaction monitoring (PRM) was evaluated by analysis of replicate digests using
423 an analytical variant called QuanDirect™ (QD) (Figure 1) [32]. This method differs from classical PRM by
424 triggering data acquisition using detection of the y_1 SIL amino acid instead of characterized retention
425 times. The resulting percentage of true positive identifications for Lumos-QD was 33.3% and the false
426 discovery rate was also highest, at 12.7% (Figure 3). The high false discovery rate was due to higher
427 levels of carry-over of the peptides GVALSNVIHK, GVALSNVVHK, and DLNMDCMVAEIK, including in the
428 blanks. Average GVP detection from the 24-GVP panel was 9.5 ± 0.7 and detection sensitivity
429 (TP/(TP+FN)) was 54.5% (Figure 3 and S3). Overall, the methodologies we applied to the Lumos-QD
430 platform revealed a 1.3-fold increase in detection sensitivity in comparison to the traditional
431 methodologies we applied to the QE+ platform (Figure 3). However, the peptides
432 ARPLEQAVAAIVCTFQEYAGR and AKPLEQAVAAIVCTFQEYAGR were both detected inconsistently in the
433 Lumos-QD series when compared to the three other methods.

434 Overall, Lumos-QD peaks appear sharp and symmetrical with peak intensities that averaged at
435 three orders of magnitude and an ion current signal to noise ratio above 1:10 (Figure 2). As expected,
436 peaks identified using Lumos-QD were less stable in retention time. Standard peptide peaks drifted
437 between runs by an average variance of 20 sec (or 0.4% of total run time), which is longer than the
438 average peak width of 15 sec (or 0.3% total run time) (Figure S2). Peak variance was on average 1.5x
439 larger for the Lumos-QD method compared to the QQQ method, described below. Peak drift between a
440 standard peptide and its corresponding endogenous peptide in each run was minimal for both PRM and
441 MRM analyses (Figure S2B/C).

442 The QuanDirect method addresses a major weakness of PRM, namely retention time variability
443 that results from low flow chromatography. The more recent SureQuant methodology [57] continues
444 this mass triggering approach by searching for the internal standard precursor ion in a fast and low-
445 resolution watch mode and switches to a high-resolution quantitative mode when the isotope-labeled

Alternative Platforms for Proteomic Genotyping

446 precursor ion is detected [57]. For QD, not having to pre-determine strict elution windows saves time
447 and effort. However, the traditional PRM methodology when it is optimized, and takes full advantages of
448 SIL characteristics and instrument cycle windows, and may be more sensitive and optimized. The PRM
449 acquisition method is easier to establish since it is not limited by prior identification of targeted
450 transitions. Of course, targeted acquisition only acquires limited information and therefore cannot be
451 used retrospectively.

452

453 *3.4 Analysis of Performance for the QQQ Platform*

454 Multiple reaction monitoring (MRM) was conducted using an Agilent 1290 Infinity series HPLC
455 system coupled to an Agilent 6495 triple quadrupole mass spectrometer (Figure 1). In this targeted
456 acquisition experiment, a panel of 24 SIL GVPs (Table 1, Figure S4) was added to provide a direct
457 comparison of transition signals and retention times for endogenous GVPs. The resulting percentage of
458 true positive identification was 42.4% and the false discovery rate ($FP/(FP+TP)$) was 4.7% (Figure 3).
459 Average GVP detection from the 24-GVP panel was 11.7 ± 1.7 and detection sensitivity ($TP/(TP+FN)$)
460 increased to 69.3% (Figure 3 and S3). Overall, the methodologies we applied to the QQQ platform
461 revealed a 1.6-fold increase in detection sensitivity in comparison to the traditional methodology we
462 applied to the QE+ platform. Peak shape was more uniform in the QQQ run where retention time drifted
463 between runs by an average variance of 1.5 sec, which is shorter than the average peak width (4.8 sec)
464 (Figure S2A). Peak drift between a standard peptide and its corresponding endogenous peptide in each
465 run was minimal (Figure S2B/C). The average signal to noise ratio for QQQ was lower and the overall
466 detected peak intensities were lower by an average of 3 orders of magnitude.

467 The QQQ method resulted in noisier peaks due to the smaller number of transitions selected,
468 lower mass accuracy, and shorter run times that may have resulted in overlap with extraneous ions.

Alternative Platforms for Proteomic Genotyping

469 However, the QQQ system provided the greatest increase in sensitivity. In terms of method
470 development, MRM on a QQQ system requires more development to deal with limited selection of
471 transition masses, detection parameters and manual optimization of acquisition parameters such as
472 collision energy, retention windows, dwell time, duty cycle, and cycle time. In terms of input material,
473 the QE+, Lumos-DIA, and Lumos-QD methods are all the same, with 1 μg of material injected due to
474 their use of nano-LC. However, the QQQ system used a volume of 10 μL , which averaged to $\sim 15 \mu\text{g}$ of
475 digested peptide material. This is a 15-fold increase in starting material, but still only 10 to 20% of a
476 protein digest from a single hair shaft (20 mm). The MRM method depends on a limited number of
477 transitions for identification. The performance of each transition therefore needs to be individually
478 evaluated and alternative transitions selected as necessary. Selection of the target GVP for a given non-
479 synonymous SNP is also a major consideration. The current approach to proteomic genotyping is based
480 on shotgun proteomics that allows genotype inference to occur from several chemical variants, or
481 'peptidoforms', of a GVP [31, 58, 59]. These result from expected but variable environmental chemical
482 modifications such as deamination, methionine oxidation and N-terminal acetylation. Selection of a
483 representative peptide will ideally occur from the peptidoform with the highest signal from samples
484 derived from a range of real-world contexts.

485 The false positives identified in this study have three potential causes. The first category is
486 genetic. This class of false positive is demonstrated by the peptide in K32 protein containing the SAP
487 T395M. An uncommon variant in K40 (W390R) results in the same genetically variant peptide sequence,
488 which was positively identified in subject E2. As described earlier, the second class of false positive
489 detection is due to instrument carry-over. SerpinB5 and K83 found in Lumos-QD may exemplify this
490 since these were also found in the method blanks and not in reagent blanks. These most likely reflect
491 instrument carry-over and not reagent contamination due to the low peptide abundance. The
492 associated peaks are smaller than that observed in other true positive hair shaft digests by more than

Alternative Platforms for Proteomic Genotyping

493 two orders of magnitude (data not shown). The degree of carry-over for any peptide marker can be
494 factored into appropriate thresholds during the development process for designating a positive
495 detection of endogenous GVP [24–26]. Assessment of inter-sample blanks is a crucial step which must
496 be included. Caution should be used to avoid a third category of false positive detection involving data
497 interpretation.

498 The presence of false positive assignments, or potential assignments, raises the issue of peptide
499 validation and what constitutes a positive determination. In targeted proteomics, positive determination
500 is more straightforward. If the retention time, precursor ion ratios, product ion ratios, and mass errors
501 are consistent to those of the SIL standard within a certain range, then the peptide is positively assigned.
502 If one of these aspects is missing, then it is still possible to validate through the other measures. For
503 example, precursor ions are missing for many GVPs in the Lumos-QD analysis including HEXB 207I,
504 GSDMA 128L, and K39 456R. However, other measurements such as product ion ratios (data not shown)
505 and retention times (Figure S2) are consistent with the standard. Therefore, these are considered
506 positive assignments. Without the use of SIL standards, as we see with the QE+ and Lumos-DIA methods,
507 this is not a straightforward task. As a first step of validation, a library may be used to compare
508 precursor and product ions. In this analysis, we use Prosit [41] as a guideline for ion ratio comparison.

509 Traditional QQQ platforms differ significantly from research mass spectrometry platforms in
510 chromatography and ionization. The QQQ platform used here employs an analytical column with
511 dimensions of 2.1 mm x 100 mm and a flowrate of 400 μ L/min. This platform also employs an Agilent jet
512 stream ion source, which offers improved instrumental sensitivity, but is not as sensitive as nanospray
513 sources. This is in comparison to the nano-LC systems of QE+ and Lumos which used column dimensions
514 of sub-100 μ m diameter and 150-250 mm lengths and a flow rates of 200-300 nL/min. These smaller
515 diameter columns with slower flow rates offer enhanced instrumental sensitivity due to entering the
516 column in a more concentrated band, therefore lessening radial dilution [50]. Using nanoflow columns,

Alternative Platforms for Proteomic Genotyping

517 ion suppression effects and scan rate limitations are reduced, and the system is more responsive to
518 temperature changes. When considering time efficiency, analytical columns offer the advantage of 15-
519 minute proteomic runs while the nanoflow systems offer around 90 minute runs. In a non-research
520 environment where time is crucial, and especially when dealing with forensic samples, the difference of
521 six to nine samples run in 90 minutes versus one sample in 90 minutes can make a large difference in
522 time efficiency and dramatically reduce the instrumentation costs per sample. Instrument downtime
523 due to maintenance and complications is also typically lower for QQQ systems.

524 Meeting the requirements of the forensic science community is an important challenge in this
525 research. The Daubert standard requires forensic evidence to meet five major milestones including
526 testing in real-world scenarios, publication and peer review, known error rates, standards to control the
527 technique's operation, and general acceptance within the forensic science community [60]. SWGDAM
528 developmental validation guidelines for genetic studies are similar, with objectives including
529 characterization of genetic markers, species specificity, sensitivity studies, stability studies, precision and
530 accuracy, case-type samples, and population studies [61]. SWGTOX gives even further guidelines on
531 mass spectrometry standards including assessments on bias and precision, calibration models,
532 instrument carry-over, inference studies, ionization suppression and enhancement, limit of detection,
533 and limit of quantitation [24]. The proposed methods in this research meet both practical and legal
534 standards. In terms of meeting the Daubert standard, GVP analysis has also undergone testing and
535 validation studies using real-world scenarios, such as pigmentation status, body site origin and time in
536 storage, peer review, and reported error [3, 17, 27, 62, 63]. This research contributes further by
537 establishing the use of peptide standards as an additional validation option to investigators [3, 17, 19,
538 27, 38, 64, 65]. To meet the SWGDAM developmental guidelines, GVP DNA markers have been
539 characterized, species specificity is checked, sensitivity is currently being studied, stability of peptides
540 has been demonstrated, precision and accuracy are reported in proteomic datasets, and population

Alternative Platforms for Proteomic Genotyping

541 studies are currently being conducted [3, 19, 27, 38, 65]. Of the targeted mass spectrometry approaches
542 taken, the QQQ mass selection windows for primary and transition ions are broader and less selective
543 than those used in parallel reaction monitoring. However, any broadening of specificity is more than
544 compensated for by the consistency of retention time, particularly in the presence of stable isotope
545 labelled (SIL) peptides. To meet SWGTOX guidelines calibration verifications, proteomic calibration
546 models, instrument carry-over criteria are being assessed and in development, or are in place. Likewise,
547 the use of exogenous SIL peptides for inference of endogenous GVPs using transition ion and signal to
548 noise ratios can be reported and available for replication (Figures S1 and S5). Additional levels of
549 validation such as establishing limits of detection and quantification are currently under investigation.

550

551 *3.5 Extrapolation of Random Match Probability*

552 Random match probabilities for the 24 peptide panel do not exceed 1 in 1000, which is to be
553 expected of a small panel. However, random match probabilities have been reported to reach up to 1 in
554 624 million from 77 detected GVPs from a single hair shaft [17]. To model what RMP estimates could be
555 if more sensitive targeted methods were applied, inferred genotypes were modeled as a function of
556 increasing detection of GVPs. The modeled genotype frequencies of each allele were randomly selected
557 from existing GVP genotype frequencies in the European reference population of the 1000 Genomes
558 Project Consortium for 10 to 300 possible GVP detections [43]. One hundred iterations were completed
559 and minimum, median, and maximum 1/RMP estimates were plotted for increasing GVP levels (Figure
560 4). Previous data from single hair and 4mg hair digests were overlaid to validate a portion of the model
561 [17]. The model demonstrates wide variation in potential 1/RMP values and different numbers of
562 observed GVPs that reflect the stochastic nature of inferred genotypes from randomized alleles; not
563 every genotype contains an allele, and genotype frequencies can vary widely. Not all of the actual
564 overlaid 1/RMP values were within the minimum and maximum boundaries for estimated RMP, which

Alternative Platforms for Proteomic Genotyping

565 reflects a higher number of GVPs occurring within an open reading frame, and therefore were treated as
566 a single locus when processing actual GVP profiles [3]. The contingency of multiple GVPs in an open
567 reading frame were not incorporated into the model. Likewise, heterozygosity was also not
568 incorporated into the model, although the resulting product of genotype frequencies of two alleles (gf_{AB}
569 = $gf_A \times gf_B$) closely approximates and is slightly more conservative than the actual genotype frequency
570 ($gf_{AB} = 2AB$) (Figure S6). The maximum difference between the two equations was only 6.25% at an
571 allelic frequency of 0.5 (Figure S6). Expected 1/RMP values from a projected 1.6-fold increase in
572 detection sensitivity was indicated on the model as was observed in MRM. For an estimated 130 GVP
573 detections, projected values would range from an estimated maximum of 1 in 10^{18} (1 in 1 quintillion), to
574 a minimum of 1 in 10^{10} (1 in 10 billion), with a median of 1 in 10^{13} (1 in 10 trillion). For a 1.3-fold increase
575 in sensitivity, as observed using the QuanDirect variant of PRM, the roughly 100 GVP detections. This
576 was a significant improvement to current standards in proteomic genotyping and predicts that
577 individualization can routinely be obtained using a single human hair shaft. Based on 20 repeated
578 iterations there was an increase in median RMP of an order of magnitude per 8.8 ± 0.9 GVP detections.

579 Many assumptions were made in this model, which elicit broad estimates of RMP. The model
580 used does not perfectly reflect the method actually used to determine RMP from detected GVPs. GVPs
581 from the hair shaft are often clustered in the same gene product and effects of linkage disequilibrium,
582 accommodated in actual RMP estimates, were not taken into consideration. These differences may help
583 to explain deviations between modelled RMPs and actual values. Actual values of RMP shown in the
584 MCMC model are from previously published data that incorporate linkage disequilibrium into the RMP
585 calculation [17]. The panel of 24 synthetic stable isotope labeled (SIL) GVPs used in this study were
586 selected to represent a range of relative abundances from very frequently to rarely observed. This was
587 to ensure that changes in detection sensitivity would be reflected in the data. They are not a random
588 selection from the more than 400 validated GVPs identified to date and the observed 1.6-fold increase

Alternative Platforms for Proteomic Genotyping

589 in sensitivity therefore is contingent. The increase in detection sensitivity was associated with an
590 increase in false discovery rate, a scenario that is common in analytical chemistry. Since the panel of 24
591 GVPs chosen does not accurately represent the full set of 408 GVPs, a bias towards more discriminating
592 RMPs may exist. Mitigation of this effect was attempted by choosing GVPs that vary in detection
593 sensitivity, length, and detection. Heterozygosity also was not incorporated into the model, and since
594 this results in a slightly more conservative RMP estimates, this may account for some of the actual
595 1/RMP values being more discriminating than the model.

596

597 **4. Conclusion**

598 This work demonstrates the utility for alternative analytical platforms in proteomic genotyping
599 and establishes validation methods for the evaluation of inferred genotypes. Sample limitation, a lack of
600 opportunity for reproducibility, and more stringent criteria for peptide identification are all relevant
601 when interpreting data and communicating findings in a legal context. Maximizing relevant peptide
602 signals is critical. Previous proteomic optimization has occurred at the level of sample processing to
603 increase the release of detectable peptides from the hair matrix [17]. This study further optimized the
604 detection of genetically variant peptides by focusing on the analytical framework. A range of three basic
605 mass spectrometry approaches were utilized and associated to a reanalysis of GVP detection using
606 standard shotgun proteomics [17]. These approaches included data dependent acquisition (DDA),
607 systematic data independent acquisition (DIA), parallel reaction monitoring (PRM), and multiple reaction
608 monitoring (MRM). PRM and MRM methods of acquisition also included the addition of a panel of 24
609 stable isotope-labeled peptides to facilitate and validate GVP detection. While each method was
610 conducted on mass spectrometry platforms with suitable configurations for each method, additional
611 optimizations could still be conducted for each approach, particularly DIA. Nevertheless, the MRM
612 method performed best in terms of GVP detection, with an overall increase in detection sensitivity of

Alternative Platforms for Proteomic Genotyping

613 1.6x when compared to the traditional data dependent acquisition approach on a QE+ platform. This
614 platform incorporated more robust analytical column chromatography and triple quadrupole mass
615 spectrometry. In addition to increased sensitivity and a simplified analytical process, the ease of
616 explanation in a legal setting, and use of preestablished methods and accreditation standards currently
617 used in forensic toxicology should facilitate incorporation into the forensic community. In this study,
618 targeted methods applied to GVP detection enhanced the use of hair protein as a source of human
619 individualization, with a projected random match probability of 1 in 10 trillion if this method were
620 applied to all 408 currently identified GVPs. Detection of human- or fluid-identifying peptides currently
621 relies on MRM on triple quadrupole mass spectrometry platforms. An expansion of this targeted
622 approach to include GVPs has the potential to dramatically improve the accessibility of proteomic
623 genotyping, reducing costs and simplifying interpretation. Increased detection sensitivity will increase
624 the discrimination and therefore utility of resulting random match probabilities. The use of targeted
625 mass spectrometry may well place proteomic genotyping as a more accessible, quantitative, and legally
626 explainable tool.

627

628 **Funding**

629 This study was supported by the National Institute of Justice, Office of Justice Programs, U.S.
630 Department of Justice (Awards 2015-DN-BX-K065 and 2019-R2-CX-0051). The findings, statements, and
631 opinions expressed in this manuscript are those of the authors and do not necessarily represent the
632 opinions of the U.S. Department of Justice.

633

634 **Ethical Approval**

635 All procedures performed in studies involving human participants were in accordance with the ethical
636 standards of the institutional and/or national research committee and with the 1964 Helsinki

Alternative Platforms for Proteomic Genotyping

637 declaration and its later amendments or comparable ethical standards. All samples were collected
638 following the guidelines provided by the Institutional Review Board (IRB# 832726) and Institutional
639 Biosafety Committee (IBC) of the University of California, Davis, CA.

640

641 **Conflict of Interest**

642 The authors have declared no conflict of interest, with the exception of GJP who has a patent based on
643 the use of genetically variant peptides for human identification (US 8,877,455 B2, Australian Patent
644 2011229918, Canadian Patent CA 2794248, and European Patent EP11759843.3). The patent is owned
645 by Parker Proteomics LLC. Protein-Based Identification Technologies LLC (PBIT) has an exclusive license
646 to develop the intellectual property and is co-owned by Utah Valley University and GJP. This ownership
647 of PBIT and associated intellectual property does not alter policies on sharing data and materials. These
648 financial conflicts of interest are administered by the Research Integrity and Compliance Office, Office of
649 Research at the University of California, Davis to ensure compliance with University of California Policy.

650

651 **Acknowledgments**

652 The authors thank Dr. Bob Rice, Dr. Blythe Durbin-Johnson, Dr. Ben Moeller, and Dr. John Newman for
653 their advice. This publication was made possible, in part, with support from the UC Davis Genome
654 Center Bioinformatics Core Facility. The sequencing was carried out at the DNA Technologies and
655 Expression Analysis Cores at the UC Davis Genome Center, LC-MS/MS was supported by NIH Shared
656 Instrumentation Grant 1S10OD010786-01. We specifically acknowledge the assistance of Jie Li, Emily
657 Kumimoto, Siranoosh Ashtari, Vanessa K Rashbrook, and Lutz Froenicke.

658

659

660

661

662

Alternative Platforms for Proteomic Genotyping

663 **References**

664

665 [1] K. M. Legg, R. Powell, N. Reisdorph, R. Reisdorph, and P. B. Danielson, "Discovery of
666 highly specific protein markers for the identification of biological stains," *Electrophoresis*,
667 vol. 35, no. 21–22, pp. 3069–3078, 2014, doi: 10.1002/elps.201400125.

668 [2] A. Tailor, J. C. Waddington, X. Meng, and B. K. Park, "Mass Spectrometric and Functional
669 Aspects of Drug-Protein Conjugation," *Chem. Res. Toxicol.*, vol. 29, no. 12, pp. 1912–
670 1935, 2016, doi: 10.1021/acs.chemrestox.6b00147.

671 [3] G. J. Parker *et al.*, "Demonstration of protein-based human identification using the hair
672 shaft proteome," *PLoS One*, vol. 11, no. 9, pp. 1–26, 2016, doi:
673 10.1371/journal.pone.0160653.

674 [4] H. Yang, B. Zhou, M. Prinz, and D. Siegel, "Proteomic analysis of menstrual blood," *Mol.*
675 *Cell. Proteomics*, vol. 11, no. 10, pp. 1024–1035, 2012, doi: 10.1074/mcp.M112.018390.

676 [5] K. M. Legg, L. M. Labay, S. S. Aiken, and B. K. Logan, "Validation of a Fully Automated
677 Immunoaffinity Workflow for the Detection and Quantification of Insulin Analogs by LC-
678 MS-MS in Postmortem Vitreous Humor," *J. Anal. Toxicol.*, vol. 43, no. 7, pp. 505–511,
679 2019, doi: 10.1093/jat/bkz014.

680 [6] T. Buonasera *et al.*, "A comparison of proteomic, genomic, and osteological methods of
681 archaeological sex estimation," *Sci. Rep.*, vol. 10, no. 1, pp. 1–15, 2020, doi:
682 10.1038/s41598-020-68550-w.

683 [7] L. Eckhart, S. Lippens, E. Tschachler, and W. Declercq, "Cell death by cornification,"
684 *Biochim. Biophys. Acta - Mol. Cell Res.*, vol. 1833, no. 12, pp. 3471–3480, 2013, doi:
685 10.1016/j.bbamcr.2013.06.010.

686 [8] M. E. Allentoft *et al.*, "The half-life of DNA in bone: Measuring decay kinetics in 158 dated
687 fossils," *Proc. R. Soc. B Biol. Sci.*, vol. 279, no. 1748, pp. 4724–4733, 2012, doi:
688 10.1098/rspb.2012.1745.

689 [9] M. M. Holland *et al.*, "Mitochondrial DNA Sequence Analysis of Human Skeletal Remains:
690 Identification of Remains from the Vietnam War," *J. Forensic Sci.*, vol. 38, no. 3, p.
691 13439J, 1993, doi: 10.1520/jfs13439j.

692 [10] C. Ginther, L. Issel-Tarver, and M. C. King, "Identifying individuals by sequencing
693 mitochondrial DNA from teeth," *Nat. Genet.*, vol. 2, no. 2, pp. 135–138, 1992, doi:
694 10.1038/ng1092-135.

695 [11] D. Fenyő and R. C. Beavis, "A method for assessing the statistical significance of mass
696 spectrometry-based protein identifications using general scoring schemes," *Anal. Chem.*,
697 vol. 75, no. 4, pp. 768–774, 2003, doi: 10.1021/ac0258709.

698 [12] D. Fenyő, J. Eriksson, and R. Beavis, "Mass Spectrometric Protein Identification Using the
699 Global Proteome Machine," in *Methods in Molecular Biology*, no. 10, 2010, pp. 189–202.

Alternative Platforms for Proteomic Genotyping

- 700 [13] F. Meier *et al.*, “Online parallel accumulation–serial fragmentation (PASEF) with a novel
701 trapped ion mobility mass spectrometer,” *Mol. Cell. Proteomics*, vol. 17, no. 12, pp.
702 2534–2545, 2018, doi: 10.1074/mcp.TIR118.000900.
- 703 [14] M. Pla-Roca *et al.*, “Antibody colocalization microarray: A scalable technology for
704 multiplex protein analysis in complex samples,” *Mol. Cell. Proteomics*, vol. 11, no. 4, pp.
705 1–12, 2012, doi: 10.1074/mcp.M111.011460.
- 706 [15] N. Nagaraj *et al.*, “Deep proteome and transcriptome mapping of a human cancer cell
707 line,” *Mol. Syst. Biol.*, vol. 7, no. 548, pp. 1–8, 2011, doi: 10.1038/msb.2011.81.
- 708 [16] T. Geiger, A. Wehner, C. Schaab, J. Cox, and M. Mann, “Comparative proteomic analysis
709 of eleven common cell lines reveals ubiquitous but varying expression of most proteins,”
710 *Mol. Cell. Proteomics*, vol. 11, no. 3, pp. 1–11, 2012, doi: 10.1074/mcp.M111.014050.
- 711 [17] Z. C. Goecker, M. R. Salemi, N. Karim, B. S. Phinney, R. H. Rice, and G. J. Parker, “Optimal
712 processing for proteomic genotyping of single human hairs,” *Forensic Sci. Int. Genet.*, vol.
713 47, no. December 2019, p. 102314, 2020, doi: 10.1016/j.fsigen.2020.102314.
- 714 [18] K. E. Mason, P. H. Paul, F. Chu, D. S. Anex, and B. R. Hart, “Development of a Protein-
715 based Human Identification Capability from a Single Hair,” *J. Forensic Sci.*, vol. 64, no. 4,
716 pp. 1152–1159, 2019, doi: 10.1111/1556-4029.13995.
- 717 [19] Z. Zhang *et al.*, “Sensitive Method for the Confident Identification of Genetically Variant
718 Peptides in Human Hair Keratin,” *J. Forensic Sci.*, vol. 65, no. 2, pp. 406–420, 2020, doi:
719 10.1111/1556-4029.14229.
- 720 [20] M. J. MacCoss, C. C. Wu, and J. R. Yates, “Probability based validation of protein
721 identifications using a modified SEQUEST algorithm,” *Anal. Chem.*, vol. 74, no. 21, pp.
722 5593–5599, 2002, doi: 10.1021/ac025826t.
- 723 [21] J. Zhang *et al.*, “PEAKS DB: De novo sequencing assisted database search for sensitive and
724 accurate peptide identification,” *Mol. Cell. Proteomics*, vol. 11, no. 4, pp. 1–8, 2012, doi:
725 10.1074/mcp.M111.010587.
- 726 [22] J. S. Cottrell, “Protein identification using MS/MS data,” *J. Proteomics*, vol. 74, no. 10, pp.
727 1842–1851, 2011, doi: 10.1016/j.jprot.2011.05.014.
- 728 [23] H. Zhu, S. Pan, S. Gu, E. Morton Bradbury, and X. Chen, “Amino acid residue specific
729 stable isotope labeling for quantitative proteomics,” *Rapid Commun. Mass Spectrom.*,
730 vol. 16, no. 22, pp. 2115–2123, 2002, doi: 10.1002/rcm.831.
- 731 [24] Scientific Working Group for Forensic Toxicology (SWGTOX), “Scientific working group for
732 forensic toxicology (SWGTOX) standard practices for method validation in forensic
733 toxicology,” *J. Anal. Toxicol.*, vol. 37, no. 7, pp. 452–474, 2013, doi: 10.1093/jat/bkt054.
- 734 [25] European Parliament and the Council of the European Union, “96/23/EC COMMISSION
735 DECISION of 12 August 2002 implementing Council Directive 96/23/EC concerning the
736 performance of analytical methods and the interpretation of results (notified under

Alternative Platforms for Proteomic Genotyping

- 737 document number C(2002) 3044)(Text with EEA relevance) (2002/657/EC),” *Off. J. Eur.*
738 *communities*, no. L 221/8, pp. 8–36, 2002, doi: 10.1017/CBO9781107415324.004.
- 739 [26] World Anti-Doping Agency (WADA), “Identification criteria for qualitative assays
740 incorporating column chromatography and mass spectrometry,” *WADA Tech. Doc. -*
741 *TD2010IDCR*, pp. 1–9, 2010, doi: TD2010IDCR.
- 742 [27] J. A. Milan *et al.*, “Comparison of protein expression levels and proteomically-inferred
743 genotypes using human hair from different body sites,” *Forensic Sci. Int. Genet.*, vol. 41,
744 no. March, pp. 19–23, 2019, doi: 10.1016/j.fsigen.2019.03.009.
- 745 [28] R. N. Franklin, N. Karim, Z. C. Goecker, B. P. Durbin-Johnson, R. H. Rice, and G. J. Parker,
746 “Proteomic genotyping: Using mass spectrometry to infer SNP genotypes in pigmented
747 and non-pigmented hair,” *Forensic Sci. Int.*, vol. 310, 2020, doi:
748 10.1016/j.forsciint.2020.110200.
- 749 [29] S. Gallien, E. Duriez, C. Crone, M. Kellmann, T. Moehring, and B. Domon, “Targeted
750 proteomic quantification on quadrupole-orbitrap mass spectrometer,” *Mol. Cell.*
751 *Proteomics*, vol. 11, no. 12, pp. 1709–1723, 2012, doi: 10.1074/mcp.O112.019802.
- 752 [30] L. C. Gillet *et al.*, “Targeted data extraction of the MS/MS spectra generated by data-
753 independent acquisition: A new concept for consistent and accurate proteome analysis,”
754 *Mol. Cell. Proteomics*, vol. 11, no. 6, pp. 1–17, 2012, doi: 10.1074/mcp.O111.016717.
- 755 [31] C. Ludwig, L. Gillet, G. Rosenberger, S. Amon, B. C. Collins, and R. Aebersold, “Data-
756 independent acquisition-based SWATH - MS for quantitative proteomics: a tutorial,”
757 *Mol. Syst. Biol.*, vol. 14, no. 8, pp. 1–23, 2018, doi: 10.15252/msb.20178126.
- 758 [32] G. McAlister, S. Eliuk, and R. Huguet, *QuanDirect: A simplified approach to fast and*
759 *accurate, high throughput targeted MS2 quantitation using internal standards. . .*
- 760 [33] D. Remane, D. K. Wissenbach, and F. T. Peters, “Recent advances of liquid
761 chromatography–(tandem) mass spectrometry in clinical and forensic toxicology — An
762 update,” *Clin. Biochem.*, vol. 49, no. 13–14, pp. 1051–1071, 2016, doi:
763 10.1016/j.clinbiochem.2016.07.010.
- 764 [34] M. M. Mbughuni, P. J. Jannetto, and L. J. Langman, “Mass spectrometry applications for
765 toxicology,” *J. Int. Fed. Clin. Chem. Lab. Med.*, vol. 27, no. 4, pp. 2016–2043, 2016.
- 766 [35] J. Maublanc, S. Dulaurent, J. Morichon, G. Lachâtre, and J. Michel Gaulier, “Identification
767 and quantification of 35 psychotropic drugs and metabolites in hair by LC-MS/MS:
768 application in forensic toxicology,” *Int. J. Legal Med.*, vol. 129, no. 2, pp. 259–268, 2015,
769 doi: 10.1007/s00414-014-1005-1.
- 770 [36] I. Shah, A. Petroczi, M. Uvacsek, M. Ránky, and D. P. Naughton, “Hair-based rapid
771 analyses for multiple drugs in forensics and doping: Application of dynamic multiple
772 reaction monitoring with LC-MS/MS,” *Chem. Cent. J.*, vol. 8, no. 1, pp. 1–10, 2014, doi:
773 10.1186/s13065-014-0073-0.

Alternative Platforms for Proteomic Genotyping

- 774 [37] U. Garg and Y. V Zhang, "Mass Spectrometry in Clinical Laboratory: Applications in
775 Therapeutic Drug Monitoring and Toxicology," *Clin. Appl. Mass Spectrom. Drug Anal.*, vol.
776 1383, pp. 241–246, 2016, doi: 10.1007/978-1-4939-3252-8.
- 777 [38] K. F. Jones, T. L. Carlson, B. A. Eckenrode, and J. Donfack, "Assessing protein sequencing
778 in human single hair shafts of decreasing lengths," *Forensic Sci. Int. Genet.*, vol. 44, no.
779 September 2019, p. 102145, 2020, doi: 10.1016/j.fsigen.2019.102145.
- 780 [39] R. Huguet, S. Eliuk, M. Blank, V. Zabrouskov, and G. McAlister, "A simplified approach to
781 fast and accurate, high throughput targeted MS2 quantitation using internal standard,"
782 2016.
- 783 [40] B. MacLean *et al.*, "Skyline: An open source document editor for creating and analyzing
784 targeted proteomics experiments," *Bioinformatics*, vol. 26, no. 7, pp. 966–968, 2010, doi:
785 10.1093/bioinformatics/btq054.
- 786 [41] S. Gessulat *et al.*, "Prosit: proteome-wide prediction of peptide tandem mass spectra by
787 deep learning," *Nat. Methods*, vol. 16, no. 6, pp. 509–518, 2019, doi: 10.1038/s41592-
788 019-0426-7.
- 789 [42] I. W. Evett and B. Weir, *Interpreting DNA evidence: statistical genetics for forensic
790 scientists*. Sunderland MA: Sinauer Associates Sunderland MA, 1998.
- 791 [43] A. Auton *et al.*, "A global reference for human genetic variation," *Nature*, vol. 526, no.
792 7571, pp. 68–74, 2015, doi: 10.1038/nature15393.
- 793 [44] N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, "Equation
794 of State Calculations by Fast Computing Machines," *J. Chem. Phys.*, vol. 21, no. 6, pp.
795 1087–1092, 1953.
- 796 [45] W. K. Hastings, "Monte carlo sampling methods using Markov chains and their
797 applications," *Biometrika*, vol. 57, no. 1, pp. 97–109, 1970, doi: 10.1093/biomet/57.1.97.
- 798 [46] Y. Perez-Riverol *et al.*, "PRIDE inspector toolsuite: Moving toward a universal
799 visualization tool for proteomics data standard formats and quality assessment of
800 proteomexchange datasets," *Mol. Cell. Proteomics*, vol. 15, no. 1, pp. 305–317, 2016, doi:
801 10.1074/mcp.O115.050229.
- 802 [47] V. Sharma *et al.*, "Panorama: A targeted proteomics knowledge base," *J. Proteome Res.*,
803 vol. 13, no. 9, pp. 4205–4210, 2014, doi: 10.1021/pr5006636.
- 804 [48] Y. J. Lee, R. H. Rice, and Y. M. Lee, "Proteome analysis of human hair shaft: From protein
805 identification to posttranslational modification," *Mol. Cell. Proteomics*, vol. 5, no. 5, pp.
806 789–800, 2006, doi: 10.1074/mcp.M500278-MCP200.
- 807 [49] L. D. Fricker, "Limitations of Mass Spectrometry-Based Peptidomic Approaches," *J. Am.
808 Soc. Mass Spectrom.*, vol. 26, no. 12, pp. 1981–1991, 2015, doi: 10.1007/s13361-015-
809 1231-x.
- 810 [50] S. R. Wilson, T. Vehus, H. S. Berg, and E. Lundanes, "Nano-LC in proteomics: Recent

Alternative Platforms for Proteomic Genotyping

- 811 advances and approaches," *Bioanalysis*, vol. 7, no. 14, pp. 1799–1815, 2015, doi:
812 10.4155/bio.15.92.
- 813 [51] H. Henry, H. R. Sobhi, O. Scheibner, M. Bromirski, S. B. Nimkar, and B. Rochat,
814 "Comparison between a high-resolution single-stage Orbitrap and a triple quadrupole
815 mass spectrometer for quantitative analyses of drugs," *Rapid Commun. Mass Spectrom.*,
816 vol. 26, no. 5, pp. 499–509, 2012, doi: 10.1002/rcm.6121.
- 817 [52] E. V Moskovets and A. R. Ivanov, "Comparative studies of peak intensities and
818 chromatographic separation of proteolytic digests, PTMs, and intact proteins obtained
819 by nanoLC-ESI MS analysis at room and elevated temperatures," *Anal. Bioanal. Chem.*,
820 pp. 3953–3968, 2016, doi: 10.1007/s00216-016-9386-2.
- 821 [53] J. C. Gesquiere, E. Diosis, M. T. Cung, and A. Tartar, "Slow isomerization of some proline-
822 containing peptides inducing peak splitting during reversed-phase high-performance
823 liquid chromatography," *J. Chromatogr. A*, vol. 478, no. C, pp. 121–129, 1989, doi:
824 10.1016/0021-9673(89)90010-1.
- 825 [54] F. Zhang, W. Ge, G. Ruan, X. Cai, and T. Guo, "Data-Independent Acquisition Mass
826 Spectrometry-Based Proteomics and Software Tools: A Glimpse in 2020," *Proteomics*,
827 vol. 1900276, pp. 1–12, 2020, doi: 10.1002/pmic.201900276.
- 828 [55] A. Wolf-Yadlin, A. Hu, and W. S. Noble, "Technical advances in proteomics: New
829 developments in data-independent acquisition," *F1000Research*, vol. 5, no. 0, pp. 1–12,
830 2016, doi: 10.12688/f1000research.7042.1.
- 831 [56] L. K. Pino, S. C. Just, M. J. MacCoss, and B. C. Searle, "Acquiring and Analyzing Data
832 Independent Acquisition Proteomics Experiments without Spectrum Libraries," *Mol. Cell.*
833 *Proteomics*, vol. 19, no. 7, pp. 1088–1103, 2020, doi: 10.1074/mcp.P119.001913.
- 834 [57] L. Stopfer *et al.*, "High-density, targeted monitoring of tyrosine phosphorylation reveals
835 activated signaling networks in human tumors," *bioRxiv*, 2020, doi:
836 10.1101/2020.06.01.127787.
- 837 [58] L. M. Smith and N. L. Kelleher, "Proteoform: A single term describing protein
838 complexity," *Nat. Methods*, vol. 10, no. 3, pp. 186–187, 2013, doi: 10.1038/nmeth.2369.
- 839 [59] G. Rosenberger *et al.*, "Inference and quantification of peptidofoms in large sample
840 cohorts by SWATH-MS," *Nat. Biotechnol.*, vol. 35, no. 8, pp. 781–788, 2017, doi:
841 10.1038/nbt.3908.
- 842 [60] M. G. Farrell, "Daubert v Merrell Dow Pharmaceutircals, Inc: Epistemology and Legal
843 Process," *Cardozo Law Rev.*, vol. 15, pp. 2183–2217, 2014.
- 844 [61] Scientific Working Group on DNA Analysis Methods (SWGDM), "Validation Guidelines
845 for DNA Analysis Methods," no. December 2016. pp. 1–13, 2016, [Online]. Available:
846 www.swgdam.org.
- 847 [62] T. Borja *et al.*, "Proteomic genotyping of fingerprint donors with genetically variant

Alternative Platforms for Proteomic Genotyping

- 848 peptides," *Forensic Sci. Int. Genet.*, vol. 42, no. March, pp. 21–30, 2019, doi:
849 10.1016/j.fsigen.2019.05.005.
- 850 [63] G. Parker *et al.*, "Proteomic genotyping: Using mass spectrometry to infer SNP genotypes
851 in a forensic context," *Forensic Sci. Int. Genet. Suppl. Ser.*, vol. 7, no. 1, pp. 664–666,
852 2019, doi: 10.1016/j.fsigss.2019.10.130.
- 853 [64] Z. C. Goecker, "Forensic proteomics: extracting identifying information from problematic
854 evidence types," 2019.
- 855 [65] L. A. Catlin *et al.*, "Demonstration of a mitochondrial DNA-compatible workflow for
856 genetically variant peptide identification from human hair samples," *Forensic Sci. Int.*
857 *Genet.*, vol. 43, no. June, 2019, doi: 10.1016/j.fsigen.2019.102148.
- 858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891

Alternative Platforms for Proteomic Genotyping

892 **Figures**

893

Peptide Information		SNP information		
Standard Peptide Sequence	MW (average)	Gene	RSID	SAP
DSQECIL T ETEAR	1561.56	KRT39	rs17843021_G	T341M
DSQECIL M ETEAR	1591.65	KRT39	rs17843021_A	T341M
GIL I DTSR	883.92	HEXB	rs10805890_A	I207V
GIL V DTSR	869.90	HEXB	rs10805890_G	I207V
ALET V QER	954.96	GSDMA	rs7212938_G	V128L
ALET L QER	968.98	GSDMA	rs7212938_T	V128L
LEGEIN T YR	1104.10	KRT32	rs2071563_G	T395M
LEGEIN M YR	1134.20	KRT32	rs2071563_A	T395M
YISLIYTNYEAGKDDYVK	2163.30	GSTP1	rs1695_A	I105V
Y V SLIYTNYEAGKDDYVK	2149.27	GSTP1	rs1695_G	I105V
VSAMYSSSS C KLPSLSPVAR	2137.37	KRT35	rs743686_A	S36P
VSAMYSSSS P KLPSLSPVAR	2147.41	KRT35	rs743686_G	S36P
EHCSACGPL S R	1283.33	KRT39	rs7213256_C	R456Q
EHCSACGPL S QLLVK	1706.90	KRT39	rs7213256_T	R456Q
EWSTFAVGPGHCLQL N DR	2097.20	NEU2	rs2233391_A	H168N
EWSTFAVGPGHCLQL H DR	2120.24	NEU2	rs2233391_C	H168N
GVALSN V IHK	1045.16	SERPINB5	rs1455555_A	I319V
GVALSN V VHK	1031.13	SERPINB5	rs1455555_G	I319V
DLNMDC M VAEIK	1446.63	KRT83	rs2852464_C	I279M
DLNMDC I VAEIK	1428.59	KRT83	rs2852464_G	I279M
GAFLY E PCGVSTPVLSTGVLR	2233.48	KRT82	rs1732263_C	E452D
GAFLY D PCGVSTPVLSTGVLR	2219.45	KRT82	rs1732263_G	E452D
A* R PLEQAVAAIVCTFQEYAGR	2402.62	S100A3	rs36022742_C	R3K
A* K PLEQAVAAIVCTFQEYAGR	2374.60	S100A3	rs36022742_T	R3K

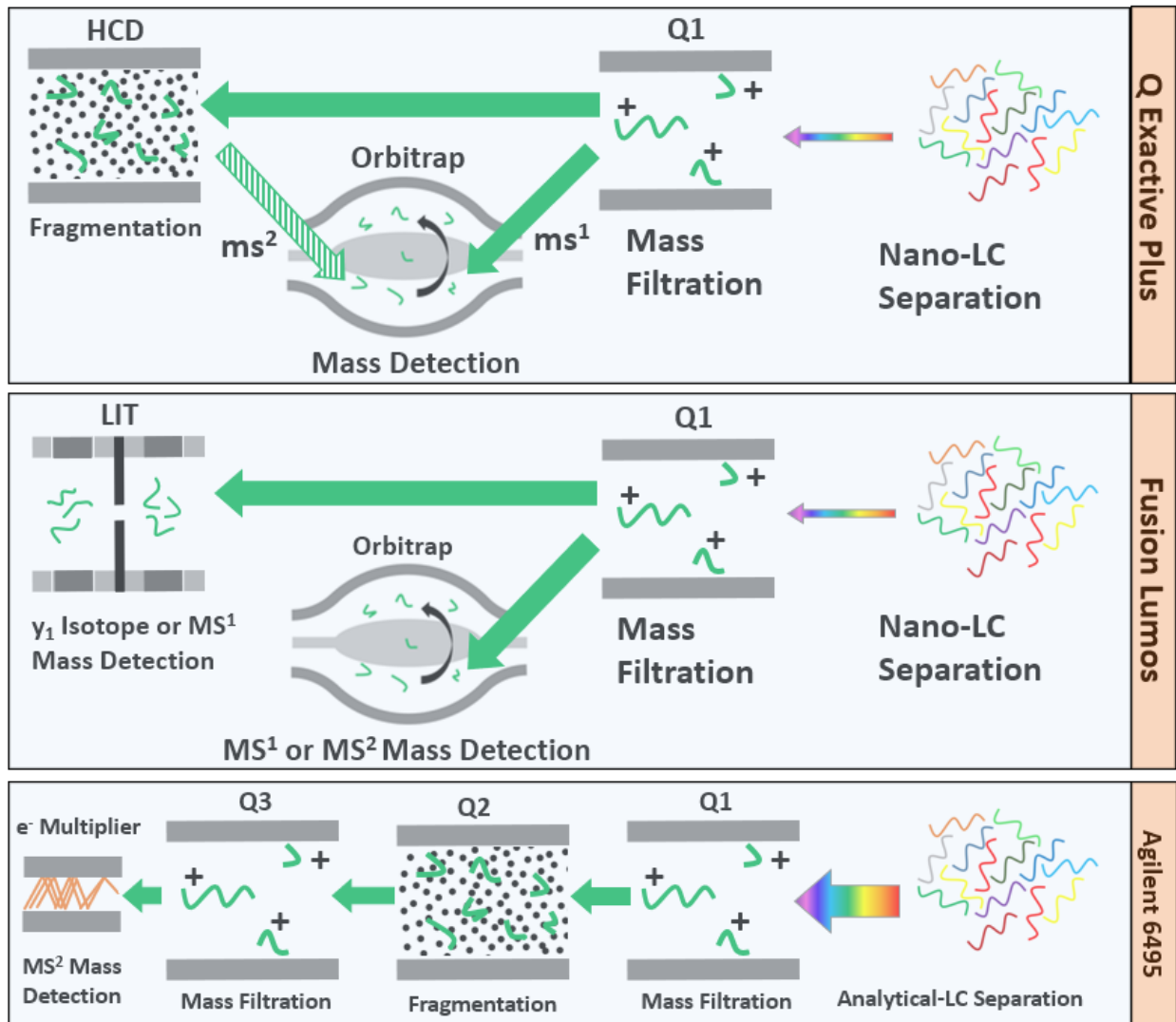
894

895 Table 1. **Genetically variant peptide standards.** These peptides were obtained from JPT Peptide
 896 Technologies and were pooled and spiked into 17 matrices from five subjects. These were spiked only
 897 into fractions being analyzed via PRM and MRM. Red amino acids indicate the SAP location per GVP, and
 898 * indicates acetylation. All cysteines (C) are carbamidomethylated (+57), and all C-terminal amino acids
 899 are isotopically labeled. R contains 6 x ¹³C and 4 x ¹⁵N (+10 Da) and K contains 6 x ¹³C and 2 x ¹⁵N (+8 Da).

900

Alternative Platforms for Proteomic Genotyping

901



902

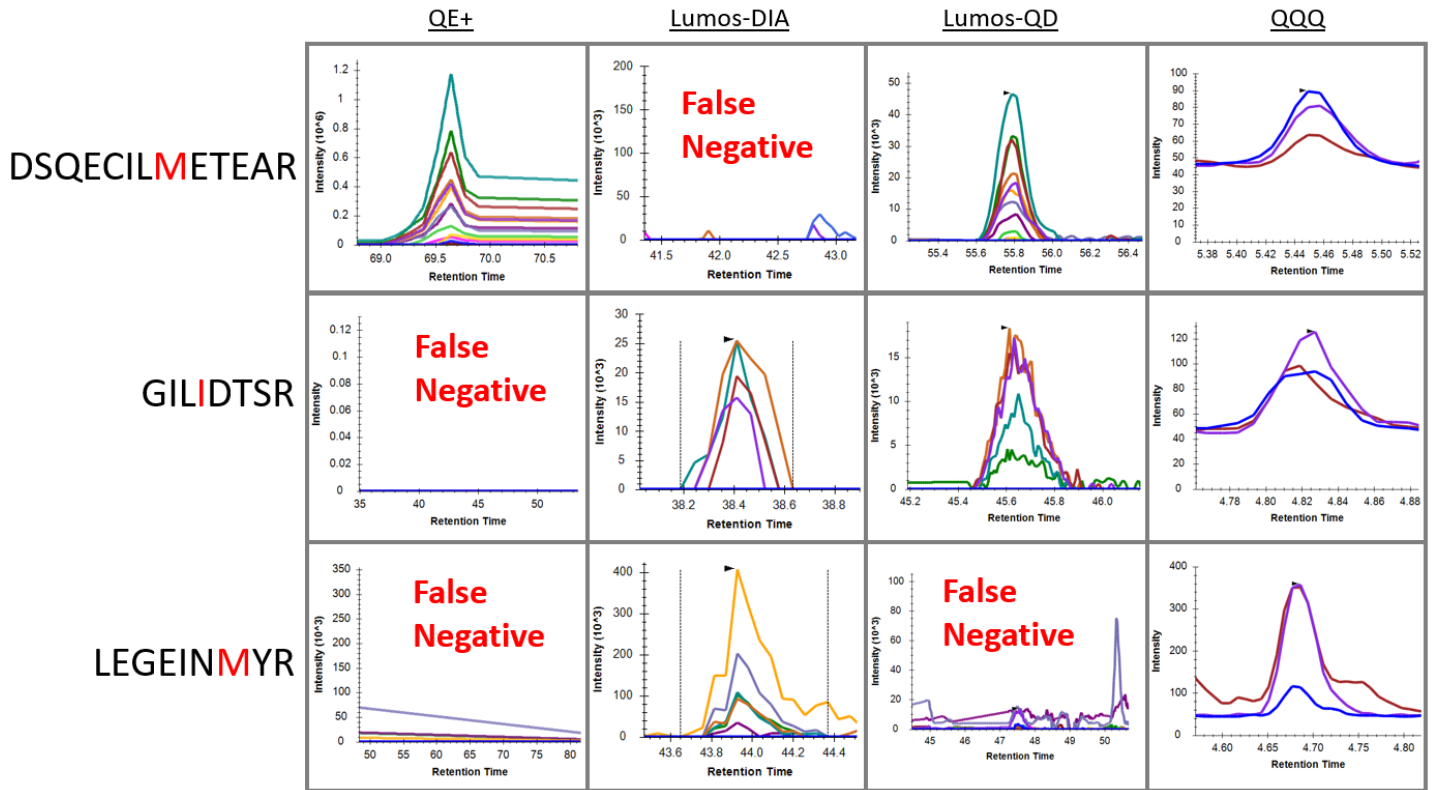
903 Figure 1. **A summary of mass spectrometry data acquisition methods.** 2 cm of scalp hair was digested.
904 This peptide mixture was then analyzed using four different LC-MS/MS data acquisition methods; two
905 methods of untargeted mass spectrometry (data-dependent acquisition on the Q Exactive plus, QE+,
906 data independent acquisition on the Fusion Lumos, Lumos-DIA) and two methods of targeted mass
907 spectrometry (parallel reaction monitoring on the Fusion Lumos, Lumos-QD, multiple reaction
908 monitoring on a triple quadrupole Agilent 6495, QQQ). For the targeted methods, an isotope-labeled
909 peptide mix was spiked into the hair digest. Both the isotope labeled peptide and endogenous peptide
910 elute together and their MS^2 spectra were compared to confirm the presence of the light isotope
911 endogenous peptide.

912

913

914

Alternative Platforms for Proteomic Genotyping



916 Figure 2. **Performance of four analytical platforms.** This figure demonstrates the usefulness of targeted
917 proteomic methods for three of the 24 GVP peptides analyzed. All peptides are expected to be present
918 in the sample as confirmed by genotyping. However, the first peptide (DSQECILMETEAR) is missing in the
919 Lumos-DIA, the second peptide (GILIDTSR) is missing in the QE+, and the third peptide (LEGEINMYR) is
920 missing in both QE+ and Lumos-QD.

921

Alternative Platforms for Proteomic Genotyping

Gene	RSID	SAP	Sequence	QE+			Lumos-DIA			Lumos-QD			QQQ		
				EUR	AFR	B	EUR	AFR	B	EUR	AFR	B	EUR	AFR	B
KRT39	rs17843021_G	T341M	DSQECILTETEAR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT39	rs17843021_A	T341M	DSQECILmETEAR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
HEXB	rs10805890_A	I207V	GILIDTSR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
HEXB	rs10805890_G	I207V	GILvDTSR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
GSDMA	rs7212938_G	V128L	ALETVQER	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
GSDMA	rs7212938_T	V128L	ALETIQER	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT32	rs2071563_G	T395M	LEGEINTYR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT32	rs2071563_A	T395M	LEGEINmYR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
GSTP1	rs1695_A	I105V	YISLIYTNYEAGKDDYVK	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
GSTP1	rs1695_G	I105V	YvSLIYTNYEAGKDDYVK	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT35	rs743686_A	S36P	VSAMYSSSCKLPSLSPVAR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT35	rs743686_G	S36P	VSAMYSSSpCKLPSLSPVAR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT39	rs7213256_C	R456Q	EHCACGPLSR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT39	rs7213256_T	R456Q	EHCACGPLSqLLVK	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
NEU2	rs2233391_A	H168N	EWSTFAVGPGHCLQLnDR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
NEU2	rs2233391_C	H168N	EWSTFAVGPGHCLQLHnDR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
SERPINB5	rs1455555_A	I319V	GVALSNVHK	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
SERPINB5	rs1455555_G	I319V	GVALSNVvHK	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT83	rs2852464_C	I279M	DLNMDCmVAEIK	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT83	rs2852464_G	I279M	DLNMDCIvAEIK	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT82	rs1732263_C	E452D	GAFLYEPCGVSTPVLSTGVLr	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
KRT82	rs1732263_G	E452D	GAFLYdPCGVSTPVLSTGVLr	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
S100A3	rs36022742_C	R3K	ARPLEQAVAAIVCTFQEYAGR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP
S100A3	rs36022742_T	R3K	AkPLEQAVAAIVCTFQEYAGR	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP	TP

QE+

26.4%	34.7%
0.7%	38.2%

FDR = 2.6% Sensitivity = 43.2%

Lumos-DIA

13.2%	47.9%
0.7%	38.2%

FDR = 5.0% Sensitivity = 21.6%

Lumos-QD

33.3%	27.8%
4.9%	34.0%

FDR = 12.7% Sensitivity = 54.5%

QQQ

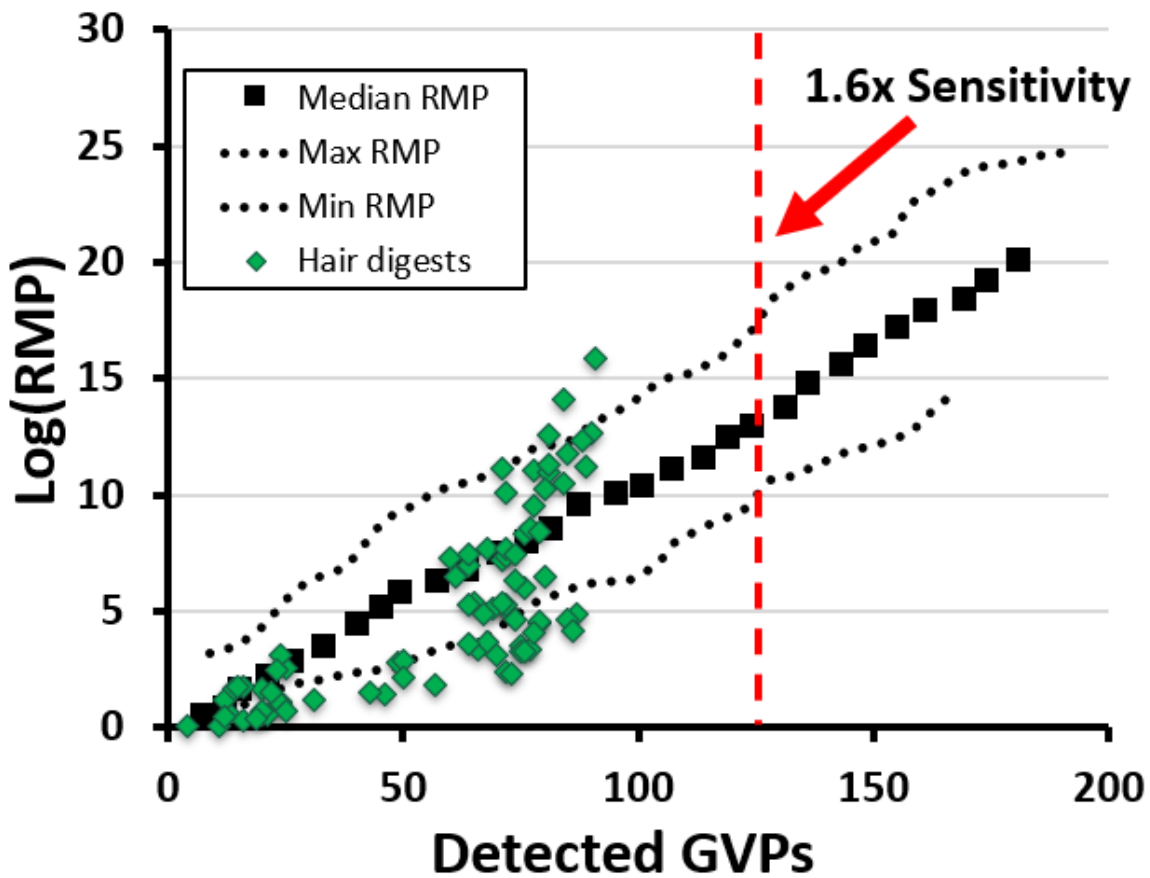
42.4%	18.8%
2.1%	36.8%

FDR = 4.7% Sensitivity = 69.3%

923 Figure 3. **GVP matrix evaluating four analytical methods.** This matrix represents GVPs that have been
 924 verified via whole exome sequencing. Each row is a variant peptide and each column is an accumulated
 925 GVP profile from three replicates. QE+, data dependent acquisition on Q Exactive+; Lumos-DIA, data
 926 independent acquisition on the Fusion Lumos; Lumos-QD, QuanDirect on Fusion Lumos; QQQ, multiple
 927 reaction monitoring on Agilent 6495; EUR, three European subjects; AFR, two African subjects; TP, true
 928 positive; FN, false negative; FP, false positive; TN, true negative; FDR, false discovery rate.

929

Alternative Platforms for Proteomic Genotyping



930

931 Figure 4. **MCMC model for RMP extrapolation.** A Markov Chain Monte Carlo model was developed to
932 estimate random match probability as a function of GVP detection. After 100 iterations, the maximum,
933 minimum, and median values were obtained. RMP values were validated from digests of 2 cm of hair
934 shaft. QQQ detection is estimated to increase GVP detection by 1.6-fold and QuanDirect™ detection is
935 estimated to increase GVP detection by 1.3-fold.

936

937

938

939

940

941

942

943