Human Forensic Identification from Fingerprint Touch Samples Using Proteogenomics



Myles W. Gardner, Ph.D.

Senior Staff Scientist, S&T Advisor Signature Science, LLC

s c i e n c e

69th Annual ASMS Conference on Mass Spectrometry and Allied Topics Forensics: Innovations and Applications, TOH am 08:30 Tuesday, November 2, 2021

Outline

- Background on Human Forensics
- IARPA Proteos Program
- General Workflow for Human ID from Fingerprints
- Sample Preparation and Chemical Analysis Methods
- Proteogenomics Data Analysis
- Results and Program Highlights
- Limitations
- Future Research Considerations



Current State of the Art – DNA Forensics



Current State of the Art – DNA Forensics



GVP = Genetically variable peptide

Challenges in Human Forensic Identification



IARPA Proteos Program

- Intelligence Advanced Research Projects Activity (IARPA) Proteos program aimed to:
 - Develop novel methods for human identification by sequencing of protein variants in touch samples
 - Develop and optimize sample preparation methods to efficiently co-extract protein and DNA in parallel
 - Augment DNA forensic identification methods when there is limited or low-quality DNA available
- 3-year program with multiple R&D and T&E phases focused on challenging developed methods against difficult sample types
 - Brass shell casings; post-IED blast material; elevator buttons; keyboard keys; wood; metal; currency; mixed contributor samples



Office of the Director of National Intelligence



https://www.iarpa.gov/index.php/research-programs/proteos

GILIDTSR ATTTTGATTGAT

GIL VDTSR



ATTTGGTTGAT

GIL V D T S R

Frequency of this Mutation in the General Population

1 in 6

ATTTTGGTTGAT



Co-Extraction of Protein and DNA from Fingerprints



DTT reduction

IAA alkylation

remove SDS

overnight

S-Trap

Trypsin digestion

Reconstitute and

LC-HRAM-MS2

analyze by

Elute peptides from

Speed-Vac to dryness

Transfer to S-Trap to

Proteomic Analysis Method | LC-HRAM-MS2

Nano-LC

- Thermo Scientific Ultimate 3000 RSLCNano
- Loading Column: Acclaim PepMap 100 C18, 75 μm × 20 mm
- Nano Column: EasySpray C18, 2.0 μm, 75 μm × 250 mm
- Mobile Phase A: 98% H₂O, 2% ACN, 0.5% FA
- Mobile Phase B: 2% H₂O, 98% ACN, 0.5% FA
- Flow Rate: 300 nL/min
- LC Run Time: 320 min (5 h gradient)
- High Resolution, Accurate MS2
 - Thermo Scientific Q Exactive Plus
 - Full MS Resolution = 70,000
 - Full MS Scan Range = m/z 375–1,575
 - Data Dependent MS2 or PRM Resolution = 17,500
 - MS2 Isolation Width = 1.6 m/z
 - PRM Inclusion List Precursor Ions = 629 (including internal standards)



Peptide-Protein Database

- Exomes translated in silico to proteomes for 50+ donors
- Compiled database of reference and variant protein sequences
- Performed *in silico* tryptic digest
- Calculated peptide frequencies
- Predicted peptide RT and MS2
- SQL database contains
 - Protein sequences
 - Peptides (frequencies)
 - Protein-peptide associations
 - Peptide-donor associations (ground truth)
 - Modified peptides (predicted RT)
 - Predicted peptide MS2
 - Genomic annotations
- Output FASTA files from database

PROTEIN_SEQUENCES	PEPTIDES_PROTEINS	PEPTIDES	E PEPTIDE_MOD_RETENTION_TIME	PRECURSOR_IONS
C A SAP CHANGE C A SAP CONC A C A C A C A C A C A C A C A C A C A	123 PEPTIDE_ID 123 PROTEIN_ID	RBC PEPTIDE_SEQ RBC GVP_HASH 123 PEPTIDE_ID 123 INTERNAL_COUNT 123 KGHGDP_COUNT 123 LOG10CHISQPVAL 123 INTERNAL_FREQ 123 RGHGDP_FREQ 123 PEPTIDE_FREQUENCY 123 RETAIN_PEPTIDE 123 LOG10CHISQPVAL_ORIGINAL 123 PEPTIDE_FREQUENCY_ORIGINAL	ABC PEPTIDE_SEQ_UNIMOD ABC MODIFICATIONS 123 PEPTIDE_MASS 123 N_VARIABLE_MODS 123 PEPTIDE_MOD_ID 123 PEPTIDE_ID 123 RT_PREDICTED_SSRC 123 RT_PREDICTED_LUDE 123 RT_PREDICTED_DEEPLC	123 CHARGE 123 PEPTIDE_MOD_ID 123 PRECURSOR_MZ 123 PRECURSOR_ID 123 PRECURSOR_ID 123 INTENSITY 123 MZ nec ION_TYPE 123 FRAG_LOC nec ION nec PREDICT_METHOD 123 PRECURSOR_ID 123 PRECURSOR_ID 123 PRETIDE_LEN
	PROTEIN_SEQUENCES PROTEIN_ID AA_SEQUENCE START_INDEX PROTEIN_LEN PEPTIDE_GENO 123 PEPTIDE_ID 123 PROTEIN_ID ABC REF_PEPTIDE_ID 123 REF_PEPTIDE_ID 123 START 123 END ABC PRE_AA ABC POST_AA 123 REF_START 124 REF_START 125 REF_S	PROTEIN_SEQUENCES PROTEIN_ID AA_SEQUENCE START_INDEX PROTEIN_LEN PEPTIDE_GENOM_ANNOTATIONS 123 PEPTIDE_ID 123 PROTEIN_ID ABC REF_PEPTIDE_ID 123 START 123 END ABC PRE_AA ABC POST_AA 123 REF_START 123 REF_END ABC SAP_CHANGE ABC SAP_CHANGE ABC SAP_PEPTIDE_CHANGE ABC SAP_PEPTIDE_CHANGE ABC MUTATION ABC DONOR_ID 123 BAD_SAP	PROTEIN_SEQUENCES 3 PROTEIN_ID 22 PEPTIDE_ID 123 PEOTEIN_ID 123 PROTEIN_IEN 123 PROTEIN_ID 123 PROTEIN_IEN 123 PROTEIN_IEN 123 PROTEIN_IEN 123 PROTEIN_IEN 123 PEPTIDE_GENOM_ANNOTATIONS 123 PROTEIN_ID 123 PROTEIN_ID 123 PROTEIN_ID 123 PROTEIN_ID 123 PROTEIN_ID 123 REF_PEPTIDE_GENOM_ANNOTATIONS 123 REF_PEPTIDE_ID 123 REF_PEPTIDE_ID 123 REF_PEPTIDE_ID 123 REF_PEPTIDE_ID 123 REF_START 123 REA_SAP_CHANGE nec MUTATION nec DONOR_ID 123 RAD_SAP	Image: PEPTIDE_SEQ PEPTIDE_ID 123 PEPTIDE_ID 123 PEPTIDE_ID 123 PEPTIDE_ID 123 PEPTIDE_ID 123 PEPTIDE_ID 123 PEPTIDE_ID 123 PEPTIDE_ID 123 PEPTIDE_ID 123 PEPTIDE_GENOM_ANNOTATIONS 123 REF_PEPTIDE_ID 123 PEPTIDE_ID 123 REF_PEPTIDE_ID 123 PEPTIDE_ID 123 REF_PEPTIDE_ID 123 PEPTIDE_ID 123 REF_PEPTIDE_ID 123 REF_PEPTIDE_ID 123 REF_PEPTIDE_ID 123 REF_START 123 REF_START 123 REF_START 123 REF_START



LC-HRAM-MS2 DDA Data Analysis Pipeline for Discovery

MetaMorpheus, MSGF+, X!Tandem, Comet + MS2PIP + Percolator

- In-house scripts written in R
- External Software Requirements (all open source)
 - X!Tandem MS-GF+ MetaMorpheus Comet PepQuery
 - MS2PIP
 Percolator
 ProteoWizard
 Elude

Detection Criteria

- Good MS2 spectral similarity to predicted MS2 (> 0.75)
- Detection by minimum of 2 search engines
- Minimum of 5 b- & y-ions
- Minimum of 10% of ion current explained by b- & y-ions



LC-HRAM-MS2 DDA Data Analysis Pipeline

MetaMorpheus, MSGF+, X!Tandem, Comet + MS2PIP + Percolator

In-house scripts written in R

Peptide

Database



Peptide

Annotations

Detection Criteria

Manual Review of

Key (Rare)

GVP Detects

Good MS2 spectral similarity to

Genetically Variable Peptide (GVP) Panel



Count*	Protein Name	
18	Hornerine (expressed in epidermis)	
16	Keratin, type II cytoskeletal 2 epidermal	
15	Keratin, type II cytoskeletal 78	
15	Desmoplakin (DP)	
13	Keratin, type I cytoskeletal 14	
13	Filaggrin-2 (FLG-2)	
12	Keratin, type II cytoskeletal 6A	
12	Filaggrin	
12	Keratinocyte proline-rich protein	
12	Junction plakoglobin (Desmoplakin-3)	
11	Keratin, type II cytoskeletal 1	
	Count* 18 16 16 15 15 13 13 12 12 12 12 12 12 12 12 12 12 12 12 12	

Common SAPs	Count
I > V	14
A > T	13
E > K	13
G > S	12
R > Q	11
T > M	10
V > I	10
N > S	10

* Includes any skin protein or gene associated with the peptide (i.e., a peptide may be associated with multiple proteins and/or genes), or total peptides that are associated with the protein † Within a set of 52 total donors

LC-HRAM-MS2 PRM Data Analysis Pipeline for Targeting

OpenSwathWorkflow + MS2PIP + Percolator

- In-house scripts written in R
- External Software Requirements (all open source)
 - OpenMS PyProphet
 MS2PIP PepQuery
 - Percolator ProteoWizard

Detection Criteria

- Good MS2 spectral similarity to predicted MS2
- Detect at least 4 (of 6 selected) MS2 transitions
- High MS1 isotope correlation
- Low MS1 mass error (8 ppm)



Results | PRM | Bulk Skin Samples Using Common GVPs

Common GVP Detection Results





Likelihood Ratio | Suspect vs. Known Donor



Results | PRM | Bulk Skin Samples Using Common GVPs

Common GVP Detection Results





Sample ID

GVP_b3a95243

GVP_b416bc01

GVP b73bea41

GVP b7d138ba

GVP bcab5f14

GVP bd9ef614 -

GVP be4bba3f

GVP c00edfac

GVP_c07220be

GVP_c10c8c43

GVP_c17fa507

GVP_c66cd9d5

GVP c898f14a

GVP_c954c36f-

GVP ca196603 ·

GVP_cc0844c8

GVP_ccd0eb92

GVP cde1ba41

GVP_ceaa907b

GVP cf55601e

Likelihood Ratio | Suspect vs. Known Donor



Unfired Brass Shell Casing | Case Study

- 9mm Brass Shell Casings (Blinded):
 - Three replicates
 - Question: Who is/are the contributor(s)?



Unfired Brass Shell Casing | 3 Contributor Mixture | PRM

20210618_Proteos_QE_028 | PR01,PR02,PR14 | commonRareAllObservedLowMz

Total Number of Peptides Detected: 181 / 518 Total Number of iRT Peptides Detected: 9 / 11 RMP Estimate (All Peptides): 3.082e-11 (173) Mean Abs Ret. Time Diff. (sec): 117.9 RMSE Ret. Time Diff (sec): 244.9 AlleleFreq SS LR | Most Likely Donor(s): PR01 AlleleFreq SS LR | Log2 LR Most Likely Donor(s): 3.672e+00 AlleleFreq SS LR | Log2 LR PR01,PR02,PR14: 6.420e+00 AlleleFreq Mixture LR | Most Likely Donor(s): PR01;PR14 AlleleFreq Mixture LR | Log2 LR Most Likely Donor(s): 6.420e+00



Unfired Brass Shell Casing | Rare/Reference GVPs | PRM

ADLTGISPSPSLYLSK | 2+ 20210618_Proteos_QE_028.mzML | PR01,PR02,PR14

m/z 824.9433 | 2+ | abs(ppm error): 0.41 | observed rt (sec): 8506 | delta rt (sec): -158 | num scans: 10 PEP: 3.33e-05 | QVALUE: 0.00e+00 | SCORE: 0.89 | num tr detected: 6 NIST MF: 0.501 | MS2 DotProdRev: 0.672 | Pearson R: 0.601 | Spectral Angle: 0.470 Primary Score: LD SCORE = 10.572 Accession: ENSP00000269491 | Peptide SAP: 10-11N>S | SNP: 63566806A>G | Peptide Freg: 6.098e-04 Detection: TRUE | Binary Class: TP









ADLTGISPSPNLYLSK | 2+





	Variable	Value
	Peptide Sequence	ADLTGISPSP <mark>S</mark> LYLSK
	Reference Peptide	ADLTGISPSPNLYLSK
	SAP (Peptide Coords)	10-11N>S
	SNP	63566806A>G
	Peptide Frequency	6.098e-04
NC2 N	Internal Count	1 (n = 52)
month	KGHGDP Count	4 (n = 6,560)
duct	Chromosome	chr18
	Protein	ENSP00000269491.1
	Transcript	ENST00000269491.5
y10 y2^	Gene	ENSG00000166634.6 (SERPINB12)

Pro

- y7^*

- y9^1





Unfired Brass Shell Casing | 3 Contributor Mixture | PRM

GVP Allele Frequency Based Likelihood Ratio Estimate 20210618_Proteos_QE_028 | PR01,PR02,PR14



LR[suspect] = rmp[best non-suspect, expected GVPs in observed GVPs] / rmp[suspect, expected GVPs in observed GVPs]

- Top 2 single-source contributors are known contributors
- Only able to identify PR01 and PR14 as all detected GVPs are consistent with this 2-person mixture
- PR01 appears to be the primary contributor (by low frequency GVP detections)

Unfired Brass Shell Casing | Case Study

- 9mm Brass Shell Casings (Blinded)
 - Three replicates
 - Question: Who is/are the contributor(s)?
- Contradictory Results
 - Protein
 - PR01 = major contributor
 - PR14 = minor contributor
 - DNA
 - PR02 = major contributor
 - Probable mixture, but insufficient data for comparison to any reference profiles



Unfired Brass Shell Casing | DNA vs. Protein



Ground Truth:

- Samples provided were mixtures
- Three contributors: PR01, PR02, and PR14
- Contradictory results were complementary
- Observed similar trends for other 3-contributor mixtures on keyboard keys and elevator buttons

DNA vs. Protein Likelihoods for Mixed Contributor Samples



Current Challenges

- Sample collection and extraction methods not in complete alignment with most forensic labs (cannot use ProK!)
- Instrumentation and subject matter experts
- No equivalent to CODIS for protein markers
 - Protein profiles can be compared to each other or a known exome/whole genome sequence
- Minimum protein input requirements (~2 µg)
- Limited GVP panel size and detectability can limit likelihood ratio values
 - Even harder considering that some alleles may not be expressed or detectable
 - May be further complicated by kinship
- Rare GVPs
 - Extremely discriminating or just a false positive?
 - How do you validate a marker you have never encountered before?

Future Directions

- Collection and extraction method optimization
- Method transfer to operational labs and independent laboratory evaluation from beginning to end
- Validation on relevant sample types
- Collection of additional data on sample matrices with known contributor(s)
- Database development (structure, format, hosting, accessibility, content)

- More Human ID by Proteomics ASMS Presentations
 - Hair and skin sample prep. optimization NIST (ThP 291)
 - Proteomics for genotyping to estimate ancestry G. Parker (ThP 092)
 - Skin variant and reference peptide spectral library from Proteos data NIST (ThP 075)
 - Human ID from touch samples using DIA LC-MS/MS Univ. of Washington & Spectragen (FP 226)



Acknowledgements

signature





Project Co-Pl Myles Gardner, Ph.D.

Project Co-Pl Curt Hewitt, Ph.D.

THE OHIO STATE UNIVERSITY



Michael Freitas, Ph.D.





August Woerner, Ph.D.

- Signature Science
 - Danielle LeSassier, Ph.D.
 - Alan Smith
 - Megan Powals
- The Ohio State University
 - Andrew Reed, Ph.D.
 - Liwen Zhang, Ph.D.
 - Maryam Baniasad
- University of North Texas
 - Ben Crysup, Ph.D.
- Lawrence Livermore National Laboratory (LLNL) –
 Forensic Science Center (Test & Evaluation Team)
- National Institute of Standards and Technology (NIST) (T&E Team)
- Intelligence Advanced Research Projects Activity (IARPA)

This research is based upon work supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via contract number 2018-18041000003. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein.

- Kathleen Schulte, M.S.
- David Joiner, M.S.
- Benjamin Ludolph



© 2021 Signature Science, LLC