



Mega Autoso	omal STR Kits











Massively Parallel Sequencing



Areas of Interest

- STRs
- Mitochondrial DNA
- HID SNPs
- Ancestry SNPs
- Phenotype SNPs
- Mixtures
- Pharmacogenetics
- Microbial Forensics
- ...going to need a bigger boat







Amplification of Library Fragments

- Bridge Amplification
- Emulsion PCR









MPS Technologies

- Increasingly user-friendly
- Highly automated workflow
- Compatible with numerous genetic marker types
- Available data analysis options
- Accurate, reliable data





Massively Parallel Sequencing



Bioinformatics

- Unprecedented access to biological data • data acquisition
- Managing biological databanks with numerous contributors and users
 - store, organize, networks
- Extracting useful information from large and dense biological data
 - manipulate, visualize
- Assembling molecular pieces into predictive models of biological systems for *in silico* experiments
 - modeling, inference
 - scientific computing: multiprocessor, faster processors

Objectives

- Sample to answer
 - Simplicity
 - Flexibility
 - Data and information management
- What do the data mean?

STR Calls

- FASTQ/FASTA Files
- Assign reads to loci
- Loci library (file) with configurations
- Reads grouped based on sequence similarities
 - Quality assessment
- Identify alleles
 - Allele library (file)
- Report Allele candidates plus stutter, etc.
- Evaluation by analyst





Our Efforts to Address Data Analyses





7



STRait Razor 3.0

- ~660x faster allele identification
- Perform "fuzzy" (approximate) string matching of anchor sequences 1: ACCT 2: ACGT 3: AC
- Portable compiled language, C++
- · functions on all major operating systems
- Multithreading capabilities
- · exploit multi-cpu computer architectures
- Enhanced descriptions of anchor sequences
- · for proper extraction of loci that are duplicated in tandem

Variation Example at D21S11 Locus

















Search for More Resolving STRs 1000 Genomes Project (raw sequences, unsorted) STR Catalog Viewer Summary of human STR variation compiled using lobSTR software

inized preferential ion, diversity of allele ELSEVIER Research paper Potential highly pol forensic identity test



Exploratory Multiplex

73 highly heterozygous loci using MPS chemistry
15 loci previously described by Phillips (2016) and others

- 451 unrelated individuals from three U.S. populations
 - Caucasian (CAU; n=155);
 - Hispanic (HIS; n=148);
 - African American (AFA; n=148)
- Each STR locus was characterized and reviewed manually for diversity using inhouse Excel workbooks
 - Alleles characterized by length and sequence
 - Population genetics analyses (heterozygosity; Hardy-Weinberg equilibrium (HWE); linkage disequilibrium (LD); random match probabilities (RMP))

In silico Mixtures

- A <u>subset of 20 loci</u> selected for comparison to the CODIS core loci
 - The current requirement is <u>20 CODIS core loci</u> for upload into the national DNA database
 - high heterozygosity (>90%)
 - Operationally problematic loci (even if heterozygosity >90%) were excluded
- 443 U.S. population samples
 - African American, (AFA; n=140, 8 incomplete profile samples removed)
 - Caucasian, (CAU; n=155)
 - Hispanic, (HIS; n=148)

CODIS Panel Loci	Exploratory Panel Loci	BEST Panel Loci					
D2S1338	D3S2406	D3S2406					
D12S391	D2S1360	D2S1360					
D1S1656	D7S3048	D7S3048					
D21S11	D8S1132	D8S1132					
D8S1179	D11S2368	D11S2368					
vWA	D15S822	D15S822					
D3S1358	D2N2	D2N2					
D18S51	D1N10	D1N10					
FGA	D12N15	D12N15					
D19S433	D1N16	D1N19					
D13S317	D1N19	D1N21					
D5S818	D1N21	D8N23					
D16S539	D8N23	D15N26					
D22S1045	D15N26	D14N56					
D7S820	D14N56	D3N61					
D2S441	D3N61	D12S1338					
CSF1PO	D4N70	D4N70					
D10S1248	D11N52	D2S1338					
TPOX	D17N32	D1S1656					
TH01	D2N43	D11N52					
Shaded cells reflect the CODIS core loci.							







<text><list-item><list-item><list-item><list-item>







HID-Ion AmpliSeq [™] Ancestry Panel							
	Sample	Biogeographic Ancestry					
	1	European					
	3	European					
	4	Asian					
	5	European					
	6	European					
	7	Asian					
	10	European					
	13	African Americans					
	14	African admix					
	15	African admix					
	16	African					
	17	European					



	I	Iaj	ologro	up	Pr	edictor	ſ	
	Yfiler		F	orenSeq		Cor	nbined	
Res	ults Tabl	e	Re	ults Tabl	e	Resu	lts Tabl	•
Haplo- group	Fitness score	Proba- bility (%)	Haplo- group	Fitness score	Proba- bility (%)	Haplo- group	Fitness score	Proba- bility (%)
E1a	2	0.0	Ela	2	0.0	Ela	1	0.0
Elbla	4	0.0	Elbla	29	0.0	Elbla	8	0.0
E1b1b	5	0.0	Elblb	18	0.0	E1b1b	9	0.0
G1	3	0.0	G1	8	0.0	G1	4	0.0
G2a	3	0.0	G2a	11	0.0	G2a	6	0.0
G2b	1	0.0	G2b	0	0.0	G2b	0	0.0
I1	2	0.0	11	6	0.0	11	3	0.0
I2a	10	0.0	I2a	35	0.0	I2a	18	0.0
I2b	2	0.0	126	10	0.0	12b	4	0.0
J1	3	0.0	J1	11	0.0	J1	5	0.0
J2a	5	0.0	J2a	27	0.0	J2a	10	0.0
J2b	1	0.0	J2b	9	0.0	J2b	3	0.0
L	7	0.0	L	21	0.0	L	9	0.0
N	7	0.0	N	13	0.0	N	7	0.0
0	10	0.0	0	26	0.0	0	14	0.0
Q	10	0.0	Q	26	0.0	Q	14	0.0
R1a	11	0.0	Rla	31	0.0	Rla	17	0.0
R1b	29	100.0	R1b	79	100.0	R1b	40	100.0
R2	6	0.0	R2	9	0.0	R2	5	0.0
Т	5	0.0	Т	17	0.0	Т	10	0.0



mtDNA is the most successful marker



Sanger Sequencing
 Current Standard for forensic mtDNA analyses Reliable
Casework Bones
• Teeth • Hair
 Not used for mixtures However exclusions could be made Databases
• CE is current standard for most forensic marker systems

Limitations of Sequencing Technology

- · Sequencing is labor intensive
- · Analysis of results is time consuming
- Costly (prices range from \$1000 to \$3000 per mtDNA sample and only HV1 and HV2)
- · Variation in intensity of peaks
- Not quantitative- impacts mixture interpretation
- · Heteroplasmy difficulties

Length Heteroplasmy ...ccaccaaacccccccccccccccccccctt... 8 C's







- Variants distribution(middle circle; n=283)
- Mean strand coverage reverse (dark) and forward (light) strand (inner circle; n=24)
- Disproportionally low coverage observed in HVII is likely an artifact of alignment to a linear reference



Variation Across the mtGenome (n=283)

- 11,607 variants
 - · defined in relation to the rCRS
- Polymorphism density clustered in HVI/HVII
 - 2,938 of the variants (25.3%)
- ~75% of variation in coding region
- Increase the value of mtDNA

Comparison of Haplogroups and Haplotypes Generated by HVI/HVII and Genome Sequence Data HVI/HVII mtGenome





mtGenome vs HVI/HVII						
	HVI/HVII mtGenome					
Sample	HaploGroup Assignment	Quality %	HaploGroup Assignment	Quality %		
USA_TX_0028	N11a	80.3	L2a1c3	93.1		
USA_TX_0052	M73'79	95.1	L3b1a+!16124	95.1		
USA_TX_0057	G3	89.3	L3b1a7	97.2		
USA_TX_0063	N2	95.2	L3e1f	95.1		
USA_TX_0108	HV0	93.9	V2	95.4		
USA_TX_0132	R0+16189	87.7	H4a1a1a1a1	97.8		
USA_TX_0174	M33c	83.6	A2+64	90.3		
USA_TX_0175	D4e1	82.8	A2+64	91.8		
USA_TX_0257	P5	95.9	H32	92.5		
Assigned by Haplogrep and Phylotree						



Mitochondrial Genome Panel

- Multiplex short amplicon system
 - Applied Biosystems Precision ID mtDNA Whole Genome Panel
- Spans entire mitochondrial genome
 - Two multiplex panels
 - Each panel contains 81 primer pairs (plus degenerate primers)
 - Tiled, overlapping pattern
 - Amplicons are ≤ 175 bps in length























	ЛРS	ana	1 Sa	noe	er S	ear	iena	ing
1		un	* 00	¹¹ 5		equ		-111 <u>8</u>
CND4								
JIIII	++		С		т			
			fwd	NS	fwd	IVS		
chr2	50,732,108	С/Т	36	33	27	20	40.52%	TWIN A
chr2	50 732 108	C	27	50	0	0	0%	TWIN B
SNP 1			ataaag <mark>c</mark> a	CATG			Ann	A
SNP 1 2F-FOR	140606_0003	a Ata	ataaaag <mark>c</mark> a	CATG	R-FOR1406	06_0003	mhm	MMM A
SNP 1 2F-FOR 2F-FOR	140606_0003 140606_0004	0 414 ST	ATAAAAG <mark>C</mark> A A		R-FOR1406	06_0003	www.	<u>~~~~~</u>





Genetic Genealogy

Joseph James DeAngelo One of the most notorious serial killers in California



Years active 1974 - 1986

East Area Rapist (N Cal) 1976-79







SNPs vs. STRs

- SNPs are mostly bi-allelic (e.g., A or T)
 Require more SNPs for identification
 - ~ 3 SNPs = 1 STR
- SNPs have very low mutation rates compared to STRs
 Great for kinship and genealogy
- SNPs have fewer artifacts than STR data
- SNPs can provide ancestry and phenotypic (outward appearance) information
- Technology allows for typing hundreds of thousands of SNPs

SNPs for Identification



~60 'random' SNPs = Full STR Profile









Genetic Genealogy



ISOGG wiki statistics: Parent/child: 3539-3748 cMs 1st cousins: 548-1034 cMs 1st cousins 1R: 248-638 cMs 2nd cousins: 101-378 cMs 2nd cousins: 2R: 43-191 cMs 3rd cousins: 43-ca 150 cMs 3rd cousins: 1R: 11.5-99 cMs More distant cousins: 5-ca 50 cMs

What is Pharmacogenomics?

- Also referred to as "Personalized Medicine"
- Melding of classical pharmacology + human genetics
- Using a patient's genotype to optimize drug therapy and minimizing toxicity



Pharmacogenentics

- Ancient Greece, Egypt, Rome

 Favism common to Central Africa and Southwest Asia
- · 1930s Phenylthiocarbamide accident and taste blindness
- · 1950s Technology and assay development
- 1975 Debrisoquine accident and description of cytochrome p450 mono-oxygenase



Four Different Phenotypes

- Poor metabolizer (PM)
 - No active copies
- Intermediate metabolizer (IM)
 1 active copy
- Extensive metabolizer (EM)
 - 2 active copies (normal)
- Ultrarapid metabolizer (UM)
 - More than 2 active copies





Forensic Example

• Codeine

- Infant died of morphine overdose at 13 days old
 - Mother was prescribed Tylenol #3 (acetaminophen and codeine)
 - · Codeine is metabolized into morphine
 - · Mother was an ultra rapid metabolizer



http://www.ilike.com/user/co

Morphine Poisoning in a Breastfed Neonate



Cytochrome p450

- Phase I metabolism enzymes
- Increase hydrophilicity of drugs and endogenous compounds
- · Influence on drug
- Extensively studied
- Metabolizer phenotype

Diversity CYP450 Enzymes

- More than 50 enzymes
- CYP1A2, CYP2C9, CYP2C19, CYP2D6, CYP3A4, and CYP3A5 enzymes metabolize 90% of drugs
- Most are expressed in the liver, but can occur in the small intestine, lungs, placenta, and kidneys

CYP2D6

- The only constitutively expressed CYP
- Accounts for ~30% of marketed drug metabolism
- Implicated in accidental overdose and idiosyncratic drug response
 - PM, IM, EM, UM
- 2015:
 - >120 alleles
 - >18 full gene duplications

CYP2D6* alleles



CYP2D6 Genotyping

- Drug Reaction
 - Targeted massively parallel sequencing (MPS)
 - Genome wide association studies (GWAS)
 - Single nucleotide polymorphism (SNP) arrays
- No potential for discovery of new/novel variants
- MPS of whole gene provides best prediction





Other Enzymes in Opiate Pathway

- CYP2D6 ~30% marketed drugs
- UGT2B7 phase II metabolism
- ABCB1 ATP dependent transporter pglycoprotein
- OPRM1 morphine and M6G receptor
- COMT enzyme in the synaptic space





	Μ	ultigenic Affects
Gene	SNP	Enzyme Activity
	-	*1A, Wild type, considered fully functional
	rs16947, rs1135840	*2D, Normal function except when duplicated
	rs35742686, rs1135824	*3A, Nonfunctional, frameshift mutation
	rs3892097, rs28371733	*4, Nonfunctional, splicing defect
CYP2D6	-	*5, Nonfunctional, complete gene deletion
	rs5030655	*6, Nonfunctional, frameshift mutation
	rs5030656	*9, Partially functional
	rs1065852	*10, Partially functional
	rs28371706, rs16947	*17, Partially functional
OPRM1	rs1799971	Decreased
LIGT2R7	rs7439366	Increased
	rs62298861	Increased
	rs2229109	Unsure
ABCB1	rs1128503	Decreased
ADEDI	rs2032582	Decreased
	rs1045642	Decreased
	rs4633	Not independently associated with activity
	rs4818	Not independently associated with activity
COMT	rs4680	Decreased
	rs2239393	Not independently associated with activity
	rs165728	Not independently associated with activity
	rs165599	Not independently associated with activity

27

Multigenic Affects							
Gene	SNP	Enzyme Activity					
	-	*1A, Wild type, considered fully functional					
	rs16947, rs1135840	*2D, Normal function except when duplicated					
	rs35742686, rs1135824	*3A, Nonfunctional, frameshift mutation					
	rs3892097, rs28371733	*4, Nonfunctional, splicing defect					
CYP2D6	-	*5, Nonfunctional, complete gene deletion					
	rs5030655 *6, Nonfunctional, frameshift mutation						
	rs5030656 *9, Partially functional						
	rs1065852 *10, Partially functional						
	rs28371706, rs16947	*17, Partially functional					
OPRM1	rs1799971	Decreased					
UCTORT	rs7439366	Increased					
001287	rs62298861	Increased					
	rs2229109	Unsure					
ARCRI	rs1128503	Decreased					
ADCDI	rs2032582	Decreased					
	rs1045642	Decreased					
	rs4633	Not independently associated with activity					
	rs4818	Not independently associated with activity					
COMT	rs4680	Decreased					
comi	rs2239393	Not independently associated with activity					
	rs165728	Not independently associated with activity					
	rs165599	Not independently associated with activity					

Microbial Forensics

• The use of scientific means to characterize microorganisms and their products for attribution purposes of a biological terrorist attack, biocrime, hoax, or accidental release of a biological agent.

• Now expanded due to advancements in massively parallel sequencing (MPS), metagenomics, and bioinformatics











28



Expansion of Microbial Forensics

- Broader Definition
- · Today's capabilities enable greater versatility



More Nonhuman Than Human

- ~10 microbial cells for every human cell
- ~5 million genes
- Single swab
- ~10,000 bacteria/cm²
 The microbiome is a high copy number genetic marker!



Personal Microbiomes · Evidence of "personal microbiomes" has been demonstrated Potential forensic applications • 16S rRNA or WGS metagenomic sequencing · No species resolution, susceptible to stochastic effects Mainly use of unsupervised methods to demonstrate that skin . microbiome signatures from touched items associate with their donors · Few studies have utilized supervised methods for the purposes of classification · Prediction of individual identification • High accuracy (> 96%) • One time point · A method had yet to be described using supervised learning approaches with strain-level features stable over time intervals https://www.sciencedail ses/2010/03/100315161



Methods of Analysis

• Many metagenomic studies use unsupervised methods

These methods know nothing about class labels

i.e., whose metagenomic sample is whose
e.g., PCA, cluster analysis

This study uses supervised

methods.

Used for prediction
Allows models to find salient features that differ between individuals





boundaries Logistic regression



hidSkinPlex Profile and Human-Specific STR and SNP Profile

Benefits of STR/SNP/mtDNA Typing by MPS

- Many Applications
- Multiplexing
- Whole mtGenome
- Mixtures
 - · Better resolution/distinguish allele v stutter
- · Can vary thresholds based on noise
- -A not an issue
- STR data are backward compatible
- · Efficient Workflow
- Data analysis tools

Hypothesis Driven Methodologies

- Gold standard limitation is most evident in mixture interpretation
- · Substantial subjectivity
- But good sign is substantial discussion
 The real strength of the field
- Present and future issues will be in hypothesis formulation, interpretation of results, documentation and communication
 - · Education and training

Allele Drop-out and Drop-in Rates

- D: the probability of drop-out of one allele of a heterozygote ($\overline{D} = 1 D$)
 - Depends on locus and DNA quantity; from 0.0 to 0.66 have been reported
 - Can be as high as 100% in a specific case
- D₂: the homozygote drop-out probability

•
$$D_2 \approx \frac{1}{2} D^2$$
;

$$\overline{D_2} = 1 - D_2$$

C: drop in probability
Some include both stutter and contamination together

Balding and Buckleton. Interpreting low template DNA profiles. Forensic Sci Int Genet. 2009 Dec;4(1):1-10

$H_1: V + S = E$

$Pr(ABC | AB, CD, H_1) = \overline{D}\overline{D}\overline{D}\overline{D}\overline{D}\overline{C}$

- A, B, C are not dropped out $\rightarrow \overline{D}$
- D is dropped out $\rightarrow D$
- No drop-in allele $\rightarrow \overline{C}$

