

Email: becky.hill@nist.gov
Phone: 301-975-4275

Carolyn R. (Becky) Hill, Margaret C. Kline, David L. Duewer, Michael D. Coble and John M. Butler

National Institute of Standards and Technology (NIST), 100 Bureau Drive, Gaithersburg, MD 20899-8314

Poster available for download from STRBase:
http://www.cstl.nist.gov/strbase/pub_pres/Hill-ISHI2012-STRloci.pdf

The original core set of 13 Combined DNA Index System (CODIS) autosomal short tandem repeat (STR) loci were selected in November 1997 and are required by the FBI for upload of DNA profiles to the national DNA database [1,2]. As the number of profiles stored in the National DNA Index System (NDIS) continues to increase each year (>11.5 million total profiles), the likelihood of adventitious matches becomes greater [3,4]. Expanding beyond the 13 core loci is critical to reduce the potential of these types of matches occurring within the database, to increase international compatibility for data sharing, and to increase discrimination power in missing persons cases [3,4]. In November 2009, the European Union adopted five new autosomal STR loci as part of their expanded European Standard Set (ESS), including D12S391, D1S1656, D2S441, D10S1248, and D22S1045 [5,6]. All five of these loci are being considered for the expansion of the U.S. core set to provide greater capabilities for international comparisons when necessary. Also, D2S1338 and D19S443 are recommended as two new additions to the original 13 core loci because almost half of the U.S. national database already contains data for these loci. Finally, it has been suggested that the DYS391 locus be added to confirm amelogenin null alleles sometimes present in DNA profiles [3,4]. In the past few years, Promega Corporation and Life Technologies have released several new next generation STR multiplex kits that enable complete coverage of all of these additional loci. These multiplex kits have been extensively tested at National Institute of Standards and Technology (NIST), allowing the probability of identity (PI) calculations to be made with different sets of loci and population statistics, including allele frequencies and heterozygosity values for each locus, determined with our set of unrelated U.S. population samples (n=1036). With this information, it has been possible to thoroughly characterize these new STR loci beyond the original 13 CODIS core loci to determine the impact that this additional information will have on database searches. A summary of these results, including STR locus population statistics for the new STR loci, are shown in order to help assess the benefits of adding additional loci to the current 13 CODIS core loci.

References:
[1] Budowle, B., et al. (1998). CODIS and PCR-based short tandem repeat loci: law enforcement tools. *Proceedings of the Second European Symposium on Human Identification*, pp. 73-88. Madison, Wisconsin: Promega Corporation. Available at <http://www.promega.com/geneticidproc/eusymp2proc/17.pdf>.
[2] Butler, J.M. (2006). Genetics and genomics of core short tandem repeat loci used in human identity testing. *J. Forensic Sci.*, 51, 253-265.
[3] Hares, D.R. (2012a). Expanding the CODIS core loci in the United States. *Forensic Sci. Int. Genet.* 6(1):e52-4.
[4] Hares, D.R. (2012b). Addendum to expanding the CODIS core loci in the United States. *Forensic Sci. Int. Genet.* 6(5):e135.
[5] Gill, P., et al. (2006). The evolution of DNA databases-Recommendations for new European STR loci. *Forensic Sci. Int.* 156: 242-244.
[6] Gill, P., et al. (2006). New multiplexes for Europe-amendments and clarification of strategic development. *Forensic Sci. Int.* 163: 155-157.
[7] Hill, C.R., et al. (2010). Strategies for concordance testing. *Profiles in DNA (Promega)*, 13(1).
[8] Hill, C.R., et al. (2011). Concordance testing comparing STR multiplex kits with a standard data set. *Forensic Sci. Int. Genet. Suppl. Ser. 3*: e188-e189.
[9] Butler, J.M., et al. (2012). Variability of new STR loci and kits in U.S. population groups. *Profiles in DNA*, (in press).

NIST Funding: Interagency Agreement 2010-DN-R-7121 between the National Institute of Justice and NIST Office of Law Enforcement Standards.
Disclaimer: Points of view are those of the authors and do not necessarily represent the official position or policies of the US Department of Justice. Certain commercial equipment, instruments and materials are identified in order to specify experimental procedures as completely as possible. In no case does such identification imply a recommendation or endorsement by the National Institute of Standards and Technology nor does it imply that any of the materials, instruments or equipment identified are necessarily the best available for the purpose.

Introduction: STR Loci and Multiplex Kits

Additional STR Loci and New Kits

Recently, Promega and Life Technologies released new STR multiplex kits to meet the needs of the planned U.S. core loci expansion [3,4]. PowerPlex Fusion (Promega) and GlobalFiler (Life Technologies) large multiplex kits both amplify 24 loci in a single reaction. With the launch of these new kits coverage of the previous CODIS 13 loci as well as the additional required (D1S1656, D2S441, D2S1338, D10S1248, D12S391, D19S433, and DYS391) and recommended (D22S1045) loci. PP Fusion also includes Penta D and Penta E, whereas GlobalFiler also includes SE33 and a Y-indel. At NIST, all of the loci present in commercial STR kits have been extensively tested with our sample set to assess the value of different combinations of loci present in these kits as well as their relative variability in these U.S. population samples.

Commonly Used STR kits World-Wide

Locus	CODIS 13	CODIS 20	ESS 12	PP 16	PP 18D	PP ES/ESX 16	PP ES/ESX 17	PP 21	PP CS7	PP Fusion	Profiler Plus	COiler	SGM Plus	SEfiler Plus	SinoFiler	MiniFiler	Identifiler	NGM	NGM SElect	GlobalFiler	
	Required loci			Promega STR kits							Life Technologies (ABI) STR kits										
D1S1656																					
F13B																					
TPOX																					
D2S441																					
D2S1338																					
D3S1358																					
FGA																					
CSF1PO																					
D5S818																					
F13A01																					
D6S1043																					
SE33																					
D7S820																					
LPL																					
D8S1179																					
Penta C																					
D10S1248																					
TH01																					
D12S391																					
vWA																					
D13S317																					
FESFPS																					
Penta E																					
D16S539																					
D18S51																					
D19S433																					
D21S11																					
Penta D																					
D22S1045																					
Amelogenin																					
DYS391																					

29 autosomal STR loci
PowerPlex Fusion 24plex STR Kit
GlobalFiler 24plex STR Kit

Locus Characteristics (Example: D1S1656)

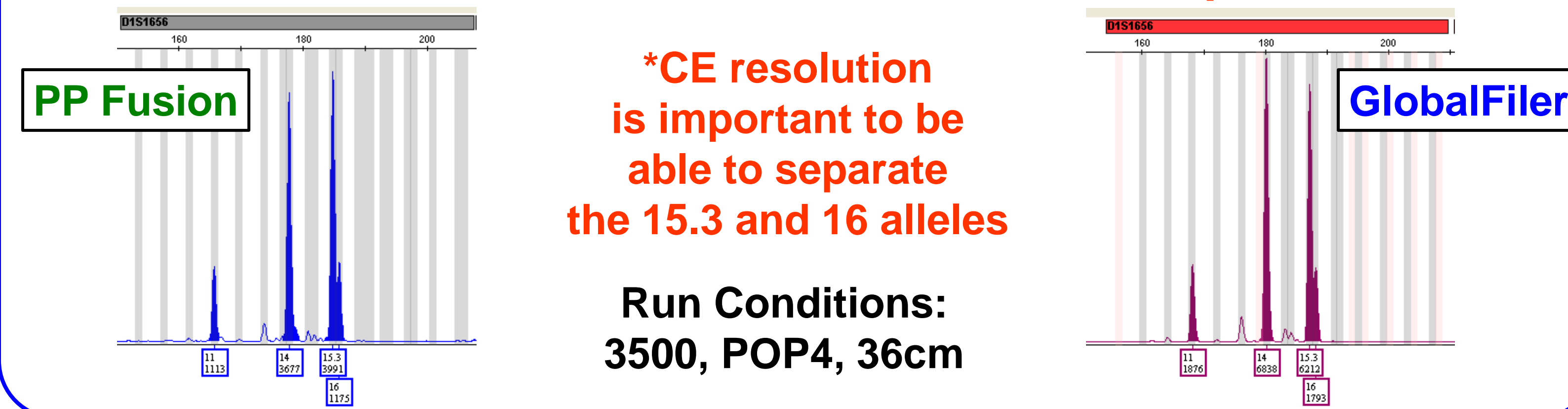
D1S1656 Allele Frequencies

STR Locus	Location	Repeat Motif	Allele Range*	# Alleles*
D2S1338	2q35	TGCC/TTCC	10 to 31	40
D19S433	19q12	AAGG/TAGG	5.2 to 20	36
Penta D	21q22.3	AAAG	1.1 to 19	50
Penta E	15q26.2	AAAG	5 to 32	53
D1S1656	1q42	TAGA	8 to 20.3	25
D12S391	12p13.2	AGAT/AGAC	13 to 27.2	52
D2S441	2p14	TCTA/TCAA	8 to 17	22
D10S1248	10q26.3	GGAA	7 to 19	13
D22S1045	22q12.3	ATT	7 to 20	14
SE33	6q14	AAAG†	3 to 49	178

Allele	African American (n=342)	Asian (n=97)	Caucasian (n=361)	Hispanic (n=236)
10	0.0146	0.0000	0.0028	0.0064
11	0.0453	0.0309	0.0776	0.0275
12	0.0643	0.0464	0.1163	0.0890
13	0.1009	0.1340	0.0665	0.1144
14	0.2573	0.0619	0.0789	0.1165
14.3	0.0073	0.0000	0.0028	0.0042
15	0.1579	0.2784	0.1496	0.1377
15.3	0.0292	0.0000	0.0582	0.0508
16	0.1096	0.2010	0.1357	0.1758
16.3	0.1023	0.0155	0.0609	0.0508
17	0.0278	0.0722	0.0471	0.0424
17.3	0.0497	0.0876	0.1330	0.1483
18	0.0029	0.0155	0.0055	0.0064
18.3	0.0234	0.0515	0.0499	0.0254
19.3	0.0073	0.0052	0.0152	0.0042

15 alleles observed

D1S1656 Mixture Profiles: 1:3 ratio, FTA spots



Materials and Methods: NIST US Population Samples

Benefits of NIST 1036 Data Set

- **Elimination of potential null alleles due to primer binding site mutations** through extensive concordance testing performed with different PCR primer sets from all available commercial STR kits
- **Ancestry testing performed** on DNA samples with autosomal SNPs, Y-SNPs, and mtDNA sequencing to verify self-declared ancestry categorization
- **Related individuals removed** based on Y-STR and mtDNA results

NIST 1036 Unrelated US Population Samples

- 1032 males + 4 females
 - 361 Caucasians (2 female)
 - 342 African Americans (1 female)
 - 236 Hispanics
 - 97 Asians (1 female)
- Anonymous donors with self-identified ancestry
 - Interstate Blood Bank (Memphis, TN) – obtained in 2002
 - Millennium Biotech, Inc. (Ft. Lauderdale, FL) – obtained in 2001
 - DNA Diagnostics Center (Fairfield, OH) – obtained in 2007
- **Complete profiles with 29 autosomal STRs + PowerPlex Y23**
 - Examined with multiple kits and in-house primer sets enabling concordance
- Additional DNA results available on subsets of these samples
 - mtDNA control region/whole genome (AFDIL)
 - >100 SNPs (AIMs), 68 InDel markers, X-STRs (AFDIL)
 - NIST assays: miniSTRs, 26plex, >100 Y-STRs, 50 Y-SNPs

Unrelated samples

All known or potential related individuals (based on autosomal & lineage marker testing) have been removed from the 1036 data set (e.g., only sons were used from father-son samples)

YSTR data available at Poster #100

Data available on STRBase: <http://www.cstl.nist.gov/biotech/strbase/NISTpop.htm>

Example of identifying and eliminating related samples

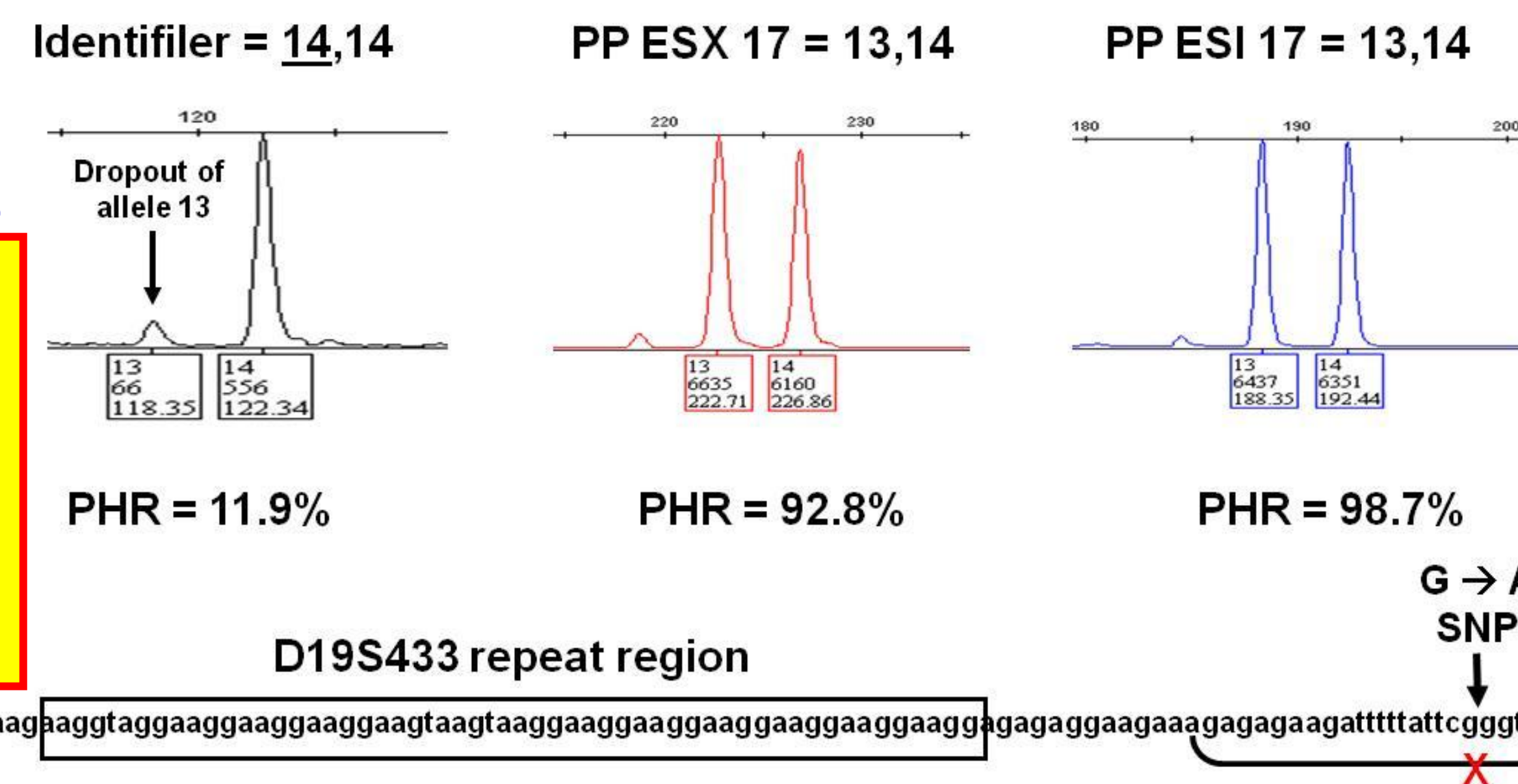
- **Hispanic samples ZT79994 and ZT79995**
 - Out of 24 autosomal STR loci, these samples share a total of 22 alleles at 22 loci (only D12S391 and Penta D have non-overlapping heterozygous alleles)
- **Full 23 Y-STR match with PowerPlex Y23**
- **Same mtDNA control region sequences**
- Kinship calculations
 - LR = 0 for parent-child
 - LR = 56,300 for full-siblings (brothers)
 - LR = 5,690 for half-siblings (or uncle-nephew, grandfather-grandson)
 - LR = 264 for first cousins
- **Decision: Remove ZT79995 from final data set**
 - ZT79994 represents this individual's family in NIST 1036

Concordance Testing with Samples [7,8]

- Many of these STR kits have different primer sequences for amplifying the same STR locus
- Need to analyze the same DNA samples with different STR typing kits looking for differences
- In some rare cases, allele dropout may occur due to mutations in primer binding regions

Current Total Concordance %

1,176,994 allele comparisons
1,225 total differences
99.90% concordance



Results and Conclusions: Characterization of STR Loci

Autosomal STR Locus Diversity with 1036 NIST Samples

Data analysis to determine individual locus diversity for each of the 29 STR loci present in commercial kits was performed with an Excel-based software tool developed by Dave Duewer at NIST to calculate allele and genotype frequencies and heterozygosities observed from the NIST 1036 data set as well as the probability of identity values reported below.

Software programs available on STRBase: <http://www.cstl.nist.gov/biotech/strbase/software.htm>

Probability of Identity [9]

- The probability of identity (P_I), also referred to as the matching probability, is **the chance that two unrelated people selected at random will have the same genotype** (first described by George Sensabaugh in 1982). The P_I value of a single locus is determined by summing the square of the observed genotype frequencies.

$$\sum_{i=1}^n x_i^2$$
 where x_i is the genotype frequency
- **Lower P_I values indicate more variability** with the genetic marker in the measured population because there are more genotypes occurring at a lower frequency.
- P_I values from independently inherited loci can be **multiplied together** to produce an expected profile P_I

STR Loci Diversity

- SE33 is the most variable locus with the highest Het_{obs} (0.9353), lowest P_I value (0.0066), and most amount of alleles and genotypes observed by over double as compared to the next highest ranked locus Penta E.
- TPOX is the least variable locus with the lowest Het_{obs} (0.06902) and highest P_I value (0.1358),
- Two of the new CODIS required loci (D2S1338 and D1S1656) rank higher than the highest ranked CODIS 13 marker (D18S51)

Loci sorted on Probability of Identity (P_I) values

Locus	Alleles Observed	Genotypes Observed	Het (obs)	P_I Value n=1036
SE33	52	304	0.9353	0.0066
Penta E	23	138	0.8996	0.0147
D2S1338	13	68	0.8793	0.0220
D1S1656	15	93	0.8890	0.0224
D18S51	22	93	0.8687	0.0258
D12S391	24	113	0.8813	0.0271
FGA	27	96	0.8745	0.0308
D6S1043	27	109	0.8494	0.0321
Penta D	16	74	0.8552	0.0382
D21S11	27	86	0.8330	0.0403
D8S1179	11	46	0.7992	0.0558
D19S433	16	78	0.8118	0.0559
vWA	11	39	0.8060	0.0611
F13A01	16	56	0.7809	0.0678
D7S820	11	32	0.7944	0.0726
D16S539	9	28	0.7761	0.0749
D13S317	8	29	0.7674	0.0765
TH01	8	24	0.7471	0.0766
Penta C	12	49	0.7732	0.0769
D2S441	15	43	0.7828	0.0841
D10S1248	12	39	0.7819	0.0845
D3S1358	11	30	0.	