

# AUTOSOMAL STR GENOMICS 101: Sequence Variation and Nomenclature

1

	<h2>AGENDA</h2>	
<p>Core Sets and Technology</p> <p>Genomic Characterization of STRs</p> <p><b>EXERCISE 1</b></p> <p>Historical Nomenclature</p> <p>STRBase and STRSeq</p> <p><b>EXERCISE 2</b></p>		<p>New Ideas in Nomenclature</p> <p>SID, STRNaming</p> <p><b>EXERCISE 3</b></p> <p>STRAND and Wrap Up</p>
		

2

# STR CORE SETS

3

## Early Development of PCR-STRs for Identification

TH01	Edwards, A., et al.	1991	Am. J. Hum. Genet. 49:746-756
TH01	Polymeropoulos, M.H., et al.	1991	Nucleic Acids Res. 19:3753
ACTBP2 (SE33)	Polymeropoulos, M.H., et al.	1992	Nucleic Acids Res. 20:1432
FGA (FIBRA)	Mills, K.A., et al.	1992	Hum. Mol. Genet. 1:779
TPOX	Anker, R., et al.	1992	Hum. Mol. Genet. 1:137
VWA (vWF)	Kimpton, C.P., et al.	1992	Hum. Mol. Genet. 1:287
D21S11	Sharma, V. and Litt, M.	1992	Hum. Mol. Genet. 1:67
Amelogenin	Sullivan, K., et al.	1993	BioTechniques 15:636-641
D3S1358	Li, H., et al.	1993	Hum. Mol. Genet. 2:1327
D18S51	Staub, R.E., et al.	1993	Genomics 15:48-56
D12S391	Lareu, M.V., et al.	1996	Gene 182:151-153



4






## Early Development of U.S. Core STR Sets

1989, 1991, 1995	TWGDAM issued guidelines for QA in DNA Analysis
1996	FBI Lab sponsors a meeting for interlaboratory validation of STR loci
1997	13 CODIS Core Loci agreed upon at STR Project Meeting
1998	NDIS implemented
1998	QAS for Forensic Testing Laboratories approved
1999	QAS for Convicted Offender DNA Databasing Laboratories approved
2011	QAS update
2015	Expanded 20 CODIS core selected
2020	QAS update



5

## Early Development of European Core STR Sets

1991	EDNAP formalizes as a WG of ISFH	
1992	EDNAP agrees on the use of STRs	
1993	FSS publishes three multiplexes, 14 STR loci (inc. vWA, TH01, D21S11)	
1995	FSS publishes SGM 6-plex: TH01, vWA, FGA, D8S1179, D18S51, and D21S11	
1996	EDNAP interlab determines TH01 and vWA would be used	
1998	Interpol establishes ESS: TH01, vWA, FGA, D21S11	
1999	Interpol expands ESS: D3S1358, D8S1179, D18S51	
2005	ENFSI and EDNAP discuss extension of ESS, propose 2 sets of 3X miniSTRs	
2006-2008	Prüm allows comparisons across European databases, adventitious matches	
2008	ENFSI meeting agreement to add: D1S1656, D2S441, D10S1248, D12S391, D22S1045	

6

# DNA Databases Worldwide

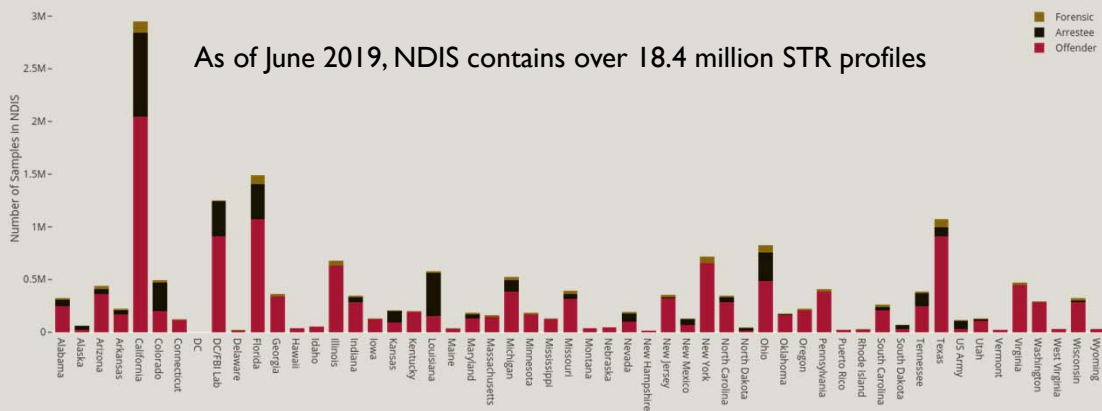
<http://dnapolicyinitiative.org/>



7

# Advances in DNA Databases

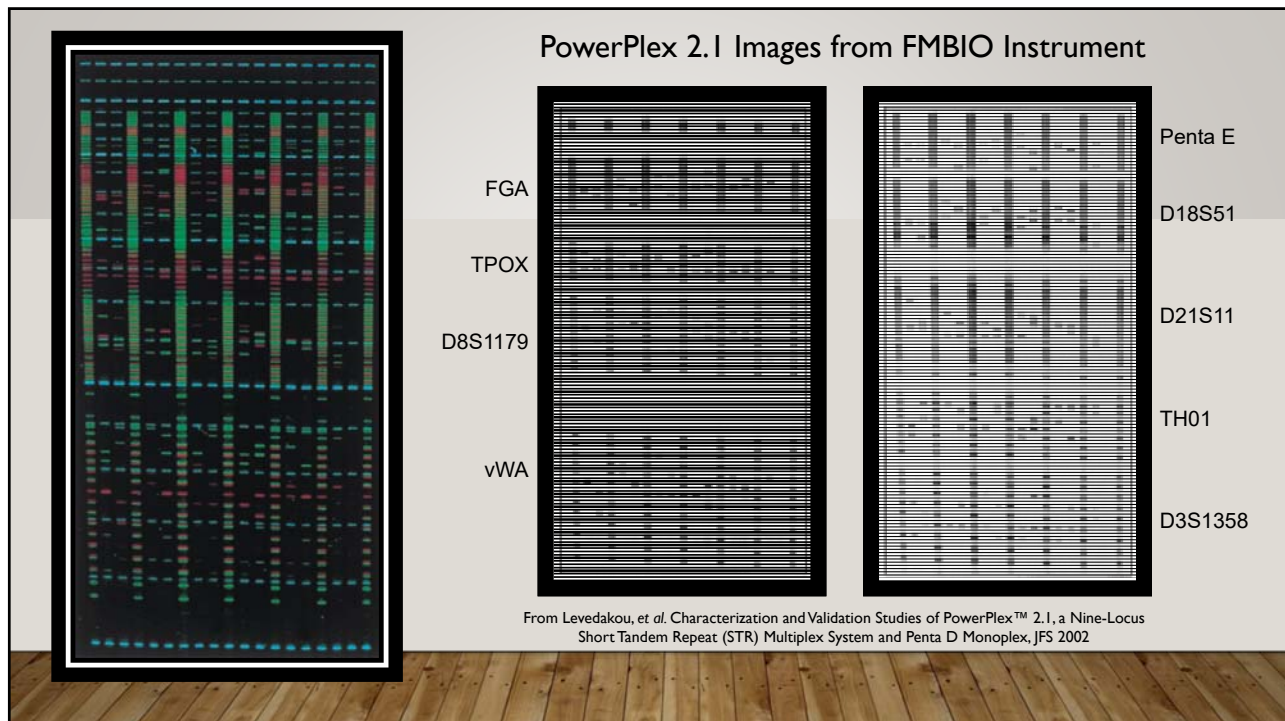
As of June 2019, NDIS contains over 18.4 million STR profiles



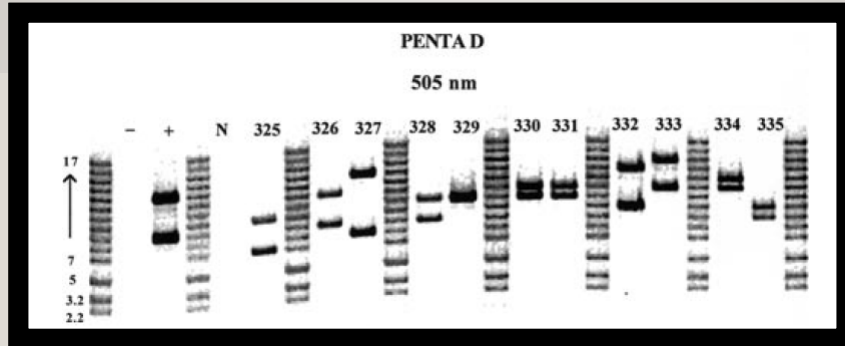
8

# STR TECHNOLOGY

9



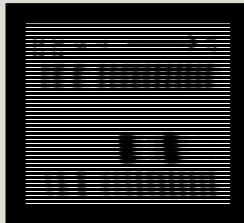
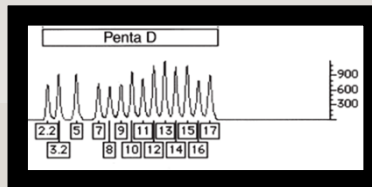
10



From Levedakou, et al. Characterization and Validation Studies of PowerPlex™ 2.1, a Nine-Locus Short Tandem Repeat (STR) Multiplex System and Penta D Monoplex, JFS 2002

11

Penta D Locus by 310  
(Capillary Electrophoresis)

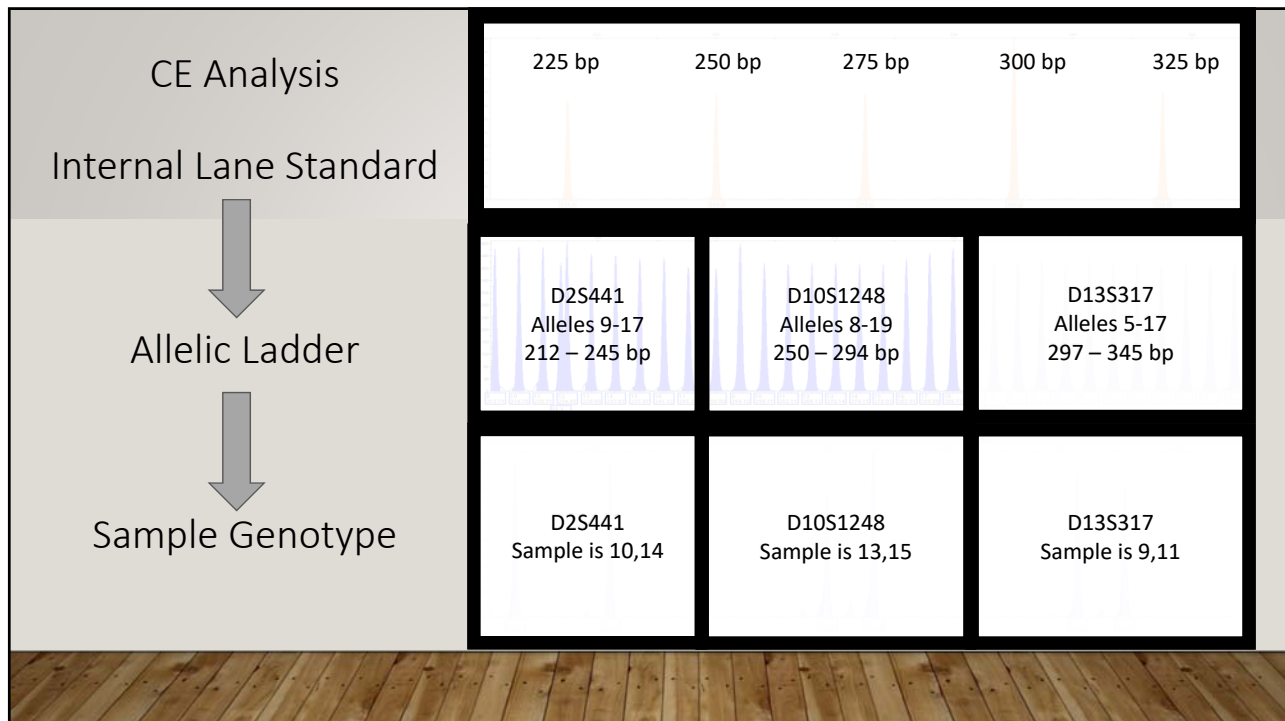


Penta D Locus by FMBIO  
(Gel Electrophoresis)

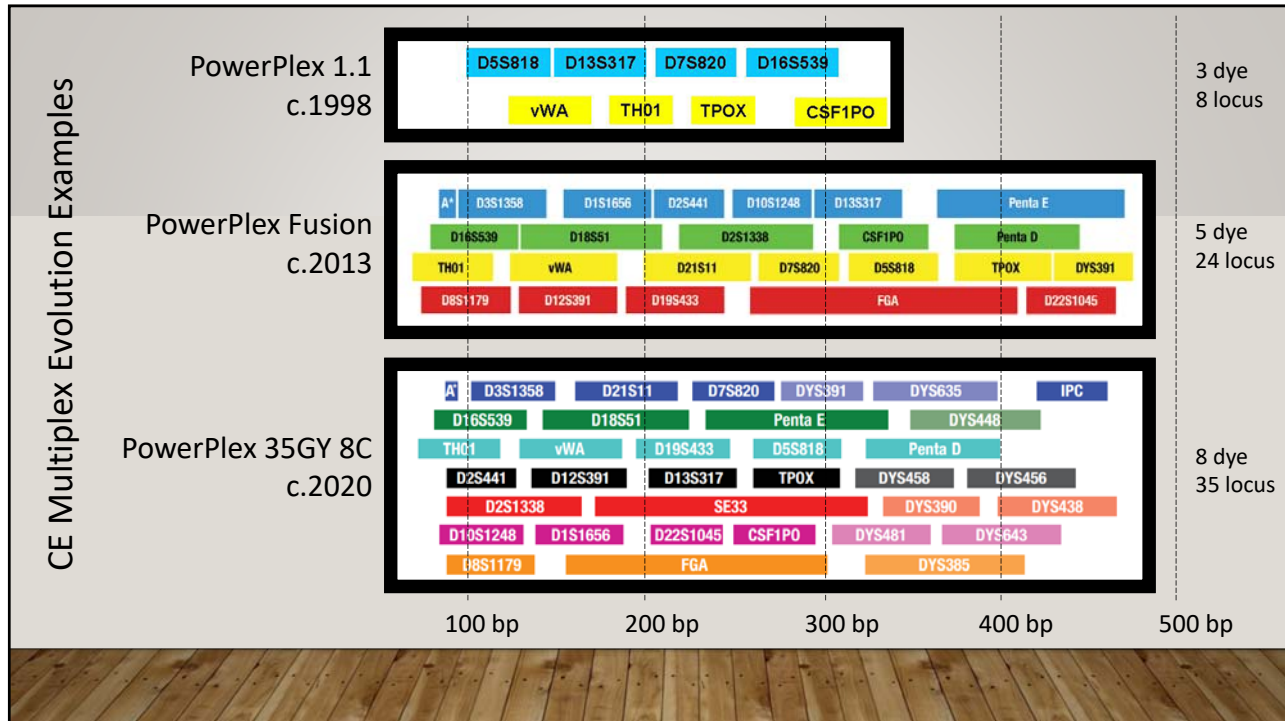
12



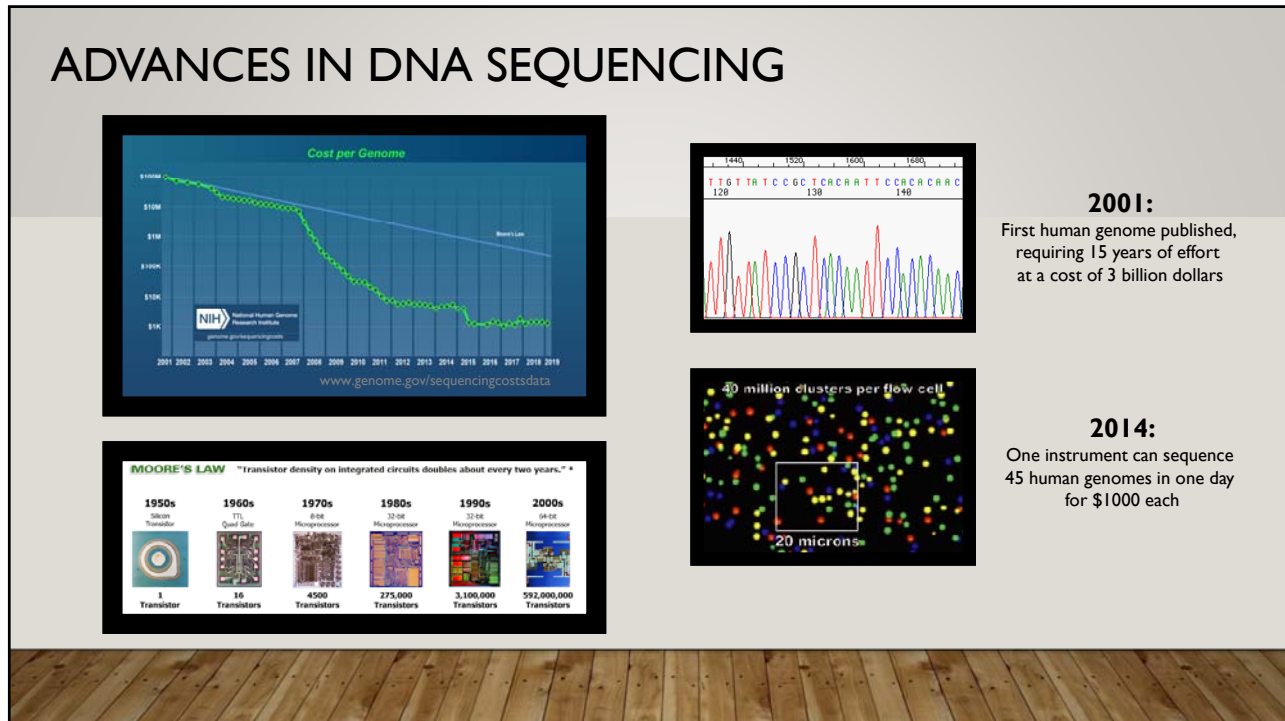
13



14



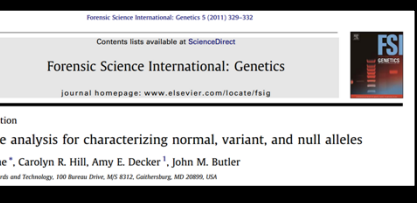
15

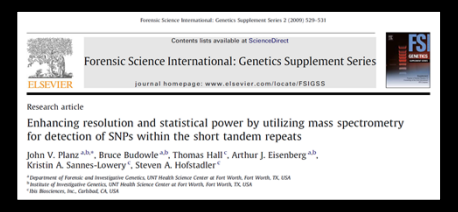


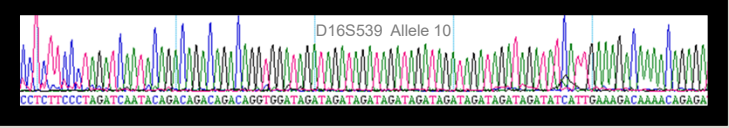
16



# SEQUENCING FORENSIC STRS







Determine the **base composition** of a PCR product containing STRs

$A_{10}G_{20}C_{12}T_4 > A_{10}G_{20}C_{11}T_5$

Provides Content not Context


17


# SEQUENCING FORENSIC STRS

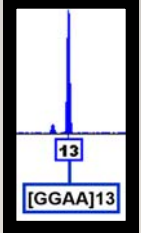
Targeted sequencing may distinguish same length alleles

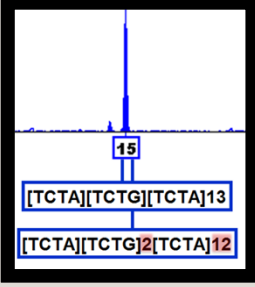
Greater degree of multiplexing

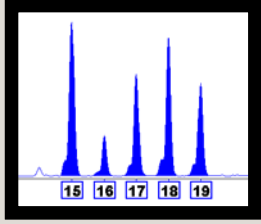
- smaller amplicons
- SNPs and other loci











18

## FORENSIC NGS KITS FOR STR TYPING



Verogen ForenSeq and FGx



Applied Biosystems  
Precision ID  
GlobalFiler NGS  
STR Panel



Promega  
PowerSeq 46GY

19

---

# GENOMIC CHARACTERIZATION OF STRs

---

20

## STR Distribution Across Human Genome



### Period Length Thresholds

Di- 11bp  $\geq$  [AT] 6  
 Tri- 14bp  $\geq$  [AAT] 5  
 Tetra- 14bp  $\geq$  [AAAT] 4  
 Penta- 16bp  $\geq$  [AAAAT] 4  
 Hexa- 17bp  $\geq$  [AAAAAT] 3

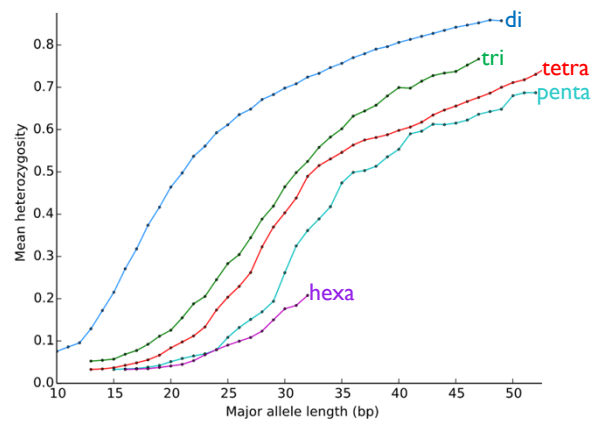
Total = 689,512 STRs

adapted from:  
 Willems T, Gymrek M, Highnam G; 1000  
 Genomes Project Consortium, Mittelman D,  
 Erlich Y. The landscape of human STR  
 variation. Genome Res. 2014 Nov; 24(11):1894-904.

21

## Heterozygosity vs Major Allele Length by Period

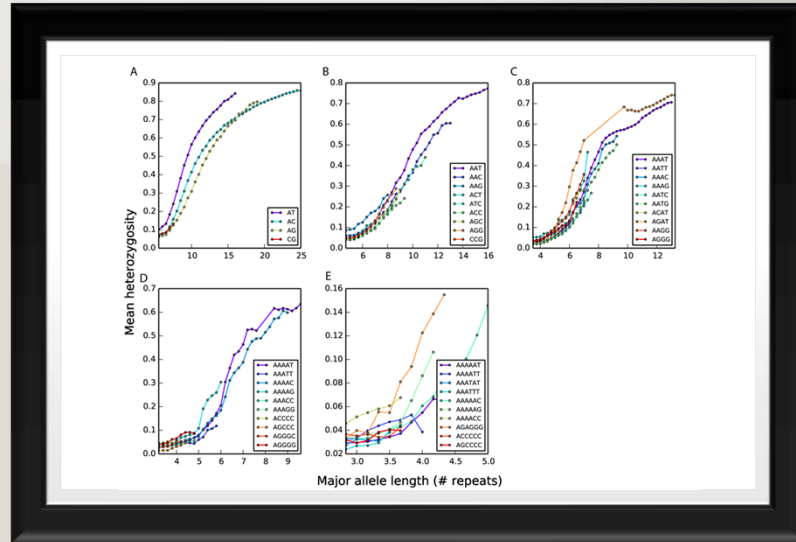
adapted from:  
 Willems T, Gymrek M, Highnam G; 1000  
 Genomes Project Consortium, Mittelman D,  
 Erlich Y. The landscape of human STR  
 variation. Genome Res. 2014 Nov;  
 24(11):1894-904.



22

## Heterozygosity vs Allele Length by Period by Motif

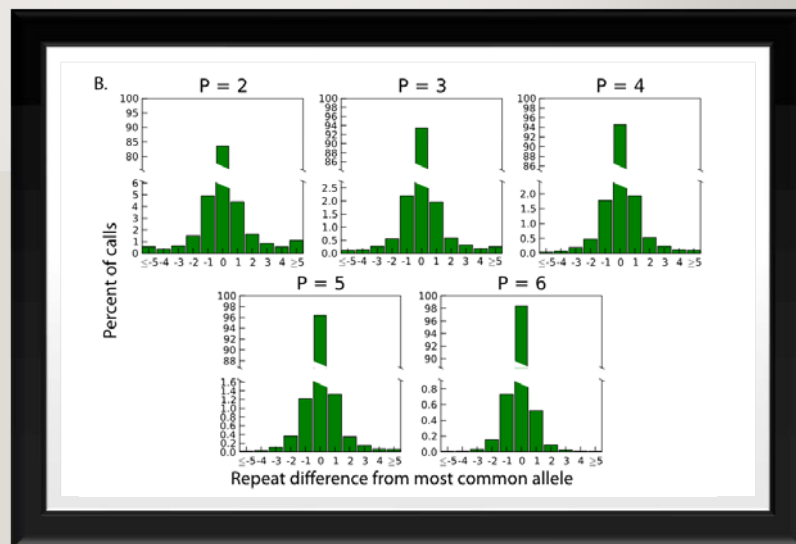
adapted from:  
 Willems T, Gymrek M, Highnam G; 1000 Genomes Project Consortium, Mittelman D, Erlich Y. The landscape of human STR variation. *Genome Res.* 2014 Nov; 24(11):1894-904.



23

## Repeat Difference from Most Common Allele by Period

adapted from:  
 Willems T, Gymrek M, Highnam G; 1000 Genomes Project Consortium, Mittelman D, Erlich Y. The landscape of human STR variation. *Genome Res.* 2014 Nov; 24(11):1894-904.



24

## REPEAT MOTIF CATEGORIES

**SIMPLE** consists of one repeat sequence

- [TCTA]<sub>n</sub>
- May have nonconsensus alleles: CCTA [TCTA]<sub>n</sub>

**COMPOUND** consists of two or more different repeat sequences

- [AGAT]<sub>n</sub> [AGAC]<sub>n</sub>
- May have nonconsensus alleles: [AGAT]<sub>n</sub> [AGAC]<sub>n</sub> AGAT

**COMPLEX** contains interspersed elements of varying period(s)

- [TCTA]<sub>n</sub> [TCTG]<sub>n</sub> [TCTA]<sub>n</sub> ta [TCTA]<sub>n</sub> tca [TCTA]<sub>2</sub> tccata [TCTA]<sub>n</sub>

25

## SIMPLE REPEATS

Consists of one repeat sequence

May have nonconsensus alleles

Locus	Allele Range	Motif	Freq	Type
D1S1656	10 to 18	CCTA [TCTA] <sub>n</sub>	54%	Simple
	10 to 17	[TCTA] <sub>n</sub>	20%	Simple
	14.3 to 19.3	CCTA [TCTA] <sub>n</sub> TCA [TCTA] <sub>n</sub>	24%	
D2S441	8 to 13	[TCTA] <sub>n</sub>	47%	Simple
	12 to 17	[TCTA] <sub>n</sub> TTTA [TCTA] <sub>2</sub>	26%	Simple
	10 to 13	[TCTA] <sub>n</sub> TCTG TCTA	18%	Simple
	11.3 to 14.3	[TCTA] <sub>n</sub> TCA [TCTA] <sub>n</sub>	5%	
D5S818	7 to 15	[ATCT] <sub>n</sub>	75%	Simple
	8 to 14	[ATCT] <sub>n</sub>	24%	Simple
	13 to 15	[ATCT] <sub>3</sub> ATGT [ATCT] <sub>n</sub>	2%	
CSF1PO	7 to 15	[ATCT] <sub>n</sub>	99%	Simple
D6S1043	8 to 15	[ATCT] <sub>n</sub>	63%	Simple
	15 to 22	[ATCT] <sub>n</sub> ATGT [ATCT] <sub>n</sub>	34%	Simple
	23 to 26	[ATCT] <sub>n</sub> ATGT [ATCT] <sub>4</sub> ATGT [ATCT] <sub>n</sub>	2%	
	18.3 to 23.3	[ATCT] <sub>n</sub> ATGT [ATCT] <sub>2</sub> ATC [ATCT] <sub>n</sub>	4%	
D7S820	6 to 14	[TATC] <sub>n</sub>	100%	Simple
D8S1179	8 to 16	[TCTA] <sub>n</sub>	38%	Simple
	11 to 17	TCTA TCTG [TCTA] <sub>n</sub>	37%	Simple
	11 to 18	[TCTA] <sub>2</sub> TCTG [TCTA] <sub>n</sub>	24%	Simple
D10S1248	8 to 19	[GGAA] <sub>n</sub>	99%	Simple
TH01	5 to 11	[AATG] <sub>n</sub>	89%	Simple
	9.3	[AATG] <sub>6</sub> ATG [AATG] <sub>3</sub>	21%	
D13S317	8 to 15	[TATC] <sub>n</sub>	99%	Simple
Penta E	5 to 25	[TCTTT] <sub>n</sub>	99%	Simple
D16S539	8 to 15	[GATA] <sub>n</sub>	100%	Simple
D18S51	9 to 24, 28	[AGAA] <sub>n</sub>	99%	Simple
	14, 15	AGAA AGCA [AGAA] <sub>n</sub>	1%	
D19S433	9 to 17	[CCTT] <sub>n</sub> CCTA CCTT CTTT CCTT	78%	Simple
	12.2 to 18.2	[CCTT] <sub>n</sub> CCTA CCTT TT CCTT	21%	
Penta D	2, 3.2, 5 to 17	[AAAGA] <sub>n</sub>	100%	Simple
D22S1045	8 to 19	[ATT] <sub>n</sub> ACT [ATT] <sub>2</sub>	100%	Simple

26

## COMPOUND REPEATS

Consists of two or more different repeat sequences

May have nonconsensus alleles

Locus	Allele Range	Motif	Freq	Type
D2S1338	15 to 24	(GGAA) <sub>n</sub> [GGCA] <sub>n</sub>	56%	Compound
	19 to 27	(GGAA) <sub>2</sub> GGAC [GGAA] <sub>n</sub> [GGCA] <sub>n</sub>	42%	Compound
	18 to 20	(GGAA) <sub>n</sub> GAAA [GGAA] <sub>2</sub> [GGCA] <sub>7</sub>	4%	Compound
D3S1358	11 to 20	TCTA [TCTG] <sub>2</sub> [TCTA] <sub>n</sub>	56%	Compound
	14 to 20	TCTA [TCTG] <sub>3</sub> [TCTA] <sub>n</sub>	28%	Compound
	12 to 18	TCTA TCTG [TCTA] <sub>n</sub>	15%	
FGA	17 to 29	(GGAA) <sub>2</sub> GGAG [AAAG] <sub>n</sub> AGAA AAAA [GAAA] <sub>3</sub>	96%	Compound
	22 to 30	(GGAA) <sub>2</sub> GGAG [AAAG] <sub>5</sub> AAGG [AAAG] <sub>n</sub> AGAA AAAA [GAAA] <sub>3</sub>	4%	Compound
	16.2 to 25.2	(GGAA) <sub>2</sub> GGAG [AAAG] <sub>n</sub> --AA AAAA [GAAA] <sub>3</sub>	2%	
vWA	11 to 21	[TAGA] <sub>n</sub> [CAGA] <sub>3-6</sub> TAGA	90%	Compound
	14 to 15	[TGGA] <sub>0,1</sub> [TAGA] <sub>3</sub> TGGA [TAGA] <sub>3</sub> [CAGA] <sub>4</sub> TAGA CAGA TAGA	10%	Compound
D12S391	14 to 27	[AGAT] <sub>n</sub> [AGAC] <sub>n</sub> AGAT	80%	Compound
	18 to 27	[AGAT] <sub>n</sub> [AGAC] <sub>n</sub>	16%	Compound
	17.3 to 19.3	AGAT GAT [AGAT] <sub>n</sub> [AGAC] <sub>7</sub> AGAT	3%	

27

## COMPLEX REPEATS

Contain interspersed elements of varying period(s) and nonconsensus alleles

Locus	Allele Range	Motif	Freq	Type
D21S11	26 to 39	[TCTA] <sub>n</sub> [TCTG] <sub>n</sub> [TCTA] <sub>n</sub> TA [TCTA] <sub>n</sub> TCA [TCTA] <sub>2</sub> TCCATA [TCTA] <sub>n</sub>	76%	Complex
	28.2 to 34.2	[TCTA] <sub>n</sub> [TCTG] <sub>n</sub> [TCTA] <sub>3</sub> TA [TCTA] <sub>n</sub> TCA [TCTA] <sub>2</sub> TCCATA [TCTA] <sub>n</sub> TA TCTA	23%	Complex
SE33	7, 11 to 23	CT [CTTT] <sub>2-3</sub> c [CTTT] <sub>n</sub> CT [CTTT] <sub>3</sub> CT [CTTT] <sub>2</sub>	47%	Complex
	19.2 to 33.2	CT [CTTT] <sub>2</sub> [CCTT] <sub>1-3</sub> C [CTTT] <sub>n</sub> TT [CTTT] <sub>n</sub> CT [CTTT] <sub>3</sub> CT [CTTT] <sub>1-2</sub>	39%	Complex

28

# CATEGORIES OF FLANKING REGIONS

1. No polymorphisms
2. Polymorphisms Associated with a Sequence Variant
3. Rare polymorphisms
4. Population or Allele Specific polymorphisms
5. "Old" polymorphisms (not population or allele specific)
6. Multiple polymorphisms in Haplotype

29

Category	Locus	rs Number	Minor Allele Frequency			Number of Associated STR Alleles			Alleles Gained with Flanking SNPs
			African American	Caucasian	Hispanic	African American	Caucasian	Hispanic	
No SNPs	D3S1358								
	FGA								
	D12S391								
	Penta E								
	D19S433								
SNPs Associated with Repeat Region Sequence Variant	D2S11								
	D151656	rs4847015	0.191	0.336	0.311	5	5	4	0
	vWA	rs11063971	0.044	0.086	0.078	1	1	1	0
		rs11063970	0.044	0.086	0.078	1	1	1	
		rs11063969	0.044	0.086	0.078	1	1	1	
rs75219269		0.044	0.086	0.078	1	1	1		
	rs199970098	0.029	-	0.011	2	-	1		
Rare SNPs	D2S441	rs74640515	0.015	0.014	0.067	1	1	2	2
	CSF1PO	C-T 36bp 5'	-	-	0.011	-	-	1	1
	D8S1179	rs138862078	0.007	-	-	1	-	-	1
	D10S1248	T-G 2bp 3'	-	0.007	-	-	1	-	1
	D18S51	rs535833682	0.029	-	-	3	-	-	1
	Penta D	rs7279663	0.022	0.014	0.011	1	2	1	2
	D22S1045	rs190864081	0.007	-	-	1	-	-	1
	Population / Allele Specific SNPs	TPOX	rs13422969	0.135	-	0.011	2	-	1
		rs115644759	0.022	-	-	2	-	-	
		rs149212737	-	0.007	-	-	1	-	
Single "Old" SNP (not population or allele specific)	TH01	rs79373318	0.148	-	-	3	-	-	3
	D16S539	rs1728369	0.404	0.171	0.278	4	4	4	4
		G-A 94bp 5'	-	-	0.011	-	-	1	5
	D2S1338	rs6736691	0.221	0.264	0.300	8	6	4	3
Multiple SNPs in Haplotype	D5S818	G-T 4bp 3'	0.316	0.257	0.189	7	5	6	17
		rs25768	0.228	0.257	0.156	5	5	5	
		rs146841551	0.029	-	-	1	-	-	
		rs541272009	0.007	-	-	1	-	-	
	D13S317	rs9546005	0.309	0.286	0.333	5	4	7	16
		rs73525369	0.066	-	-	3	-	-	
		rs73250432	-	0.014	0.011	-	2	1	
		rs146621667	0.015	-	-	2	-	-	
		4bp del 8bp 3'	-	-	0.022	-	-	2	
	D7S820	4bp del 21bp 3'	0.007	-	0.022	1	-	1	15
		rs16887642	0.177	0.050	0.022	4	2	2	
		rs7789995	0.015	0.164	0.122	2	4	4	
		rs7786079	0.176	0.021	0.011	5	2	1	
	1bp del 21bp 3'	-	-	0.011	-	-	1		

30

Category	Locus	rs Number	Minor Allele Frequency			Number of Associated STR Alleles			Alleles Gained with Flanking SNPs
			African American	Caucasian	Hispanic	African American	Caucasian	Hispanic	
No SNPs	D3S1358								
	FGA								
	D12S391								
	Penta E								
	D19S433								
	D21S11								

No SNPs were identified  
in the **PowerSeq** sequenced flanking regions of:

D3S1358                  Penta E  
FGA                        D19S433  
D12S391                 D21S11

31

Category	Locus	rs Number	Minor Allele Frequency			Number of Associated STR Alleles			Alleles Gained with Flanking SNPs
			African American	Caucasian	Hispanic	African American	Caucasian	Hispanic	
SNPs Associated with Repeat Region Sequence Variant	D1S1656	rs4847015	0.191	0.336	0.311	5	5	4	0

D1S1656 has one SNP, rs4847015

- Well distributed across populations...

**1000 Genomes Project Phase 3 allele frequencies**

Population	T (%)	C (%)
ALL	20%	80%
AFR	18%	82%
AMR	31%	69%
ASN	12%	88%
EUR	34%	66%
SAS	8%	92%

32



Category	Locus	rs Number	Minor Allele Frequency			Number of Associated STR Alleles			Alleles Gained with Flanking SNPs
			African American	Caucasian	Hispanic	African American	Caucasian	Hispanic	
SNPs Associated with Repeat Region Sequence Variant	D1S1656	rs4847015	0.191	0.336	0.311	5	5	4	0

D1S1656 has one SNP, rs4847015

- Well distributed across populations
- Well distributed across STR alleles
- Does NOT increase the number of alleles...

33

Category	Locus	rs Number	Minor Allele Frequency			Number of Associated STR Alleles			Alleles Gained with Flanking SNPs
			African American	Caucasian	Hispanic	African American	Caucasian	Hispanic	
SNPs Associated with Repeat Region Sequence Variant	D1S1656	rs4847015	0.191	0.336	0.311	5	5	4	0

D1S1656 has one SNP, rs4847015  
and it is always associated with x.3 alleles:

15.3

16.3

17.3

18.3

19.3

→

rs4847015 = T

34

Category	Locus	rs Number	Minor Allele Frequency			Number of Associated STR Alleles			Alleles Gained with Flanking SNPs
			African American	Caucasian	Hispanic	African American	Caucasian	Hispanic	
SNPs Associated with Repeat Region Sequence Variant	vWA	rs11063971	0.044	0.086	0.078	1	1	1	0
		rs11063970	0.044	0.086	0.078	1	1	1	
		rs11063969	0.044	0.086	0.078	1	1	1	
		rs75219269	0.044	0.086	0.078	1	1	1	
		rs199970098	0.029	-	0.011	2	-	1	

vWA has five SNPs

- Four are coinherited
- All four associated with a sequence variant
- Fifth SNP also associated with a sequence variant

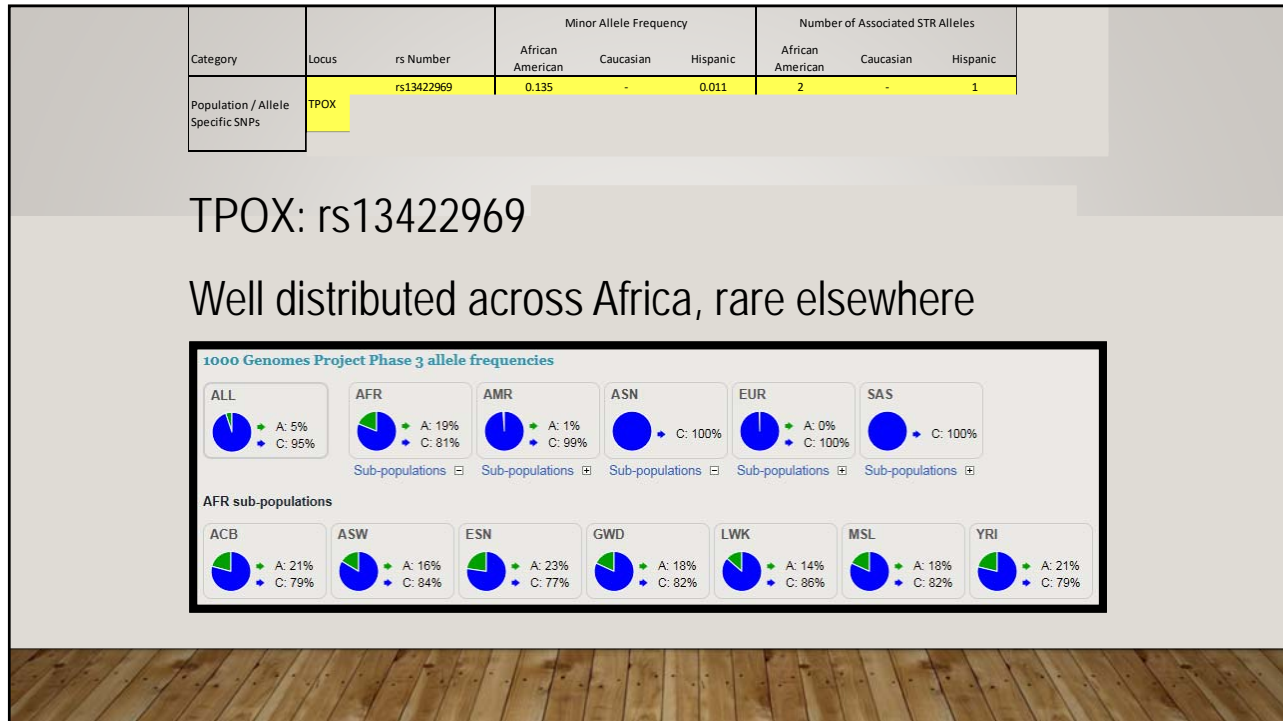
35

Category	Locus	rs Number	Minor Allele Frequency			Number of Associated STR Alleles			Alleles Gained with Flanking SNPs
			African American	Caucasian	Hispanic	African American	Caucasian	Hispanic	
Rare SNPs	D2S441	rs74640515	0.015	0.014	0.067	1	1	2	2
	CSF1PO	C-T 36bp 5'	-	-	0.011	-	-	1	1
	D8S1179	rs138862078	0.007	-	-	1	-	-	1
	D10S1248	T-G 2bp 3'	-	0.007	-	-	1	-	1
	D18S51	rs535833682	0.029	-	-	3	-	-	1
	Penta D	rs7279663	0.022	0.014	0.011	1	2	1	2
	D22S1045	rs190864081	0.007	-	-	1	-	-	1

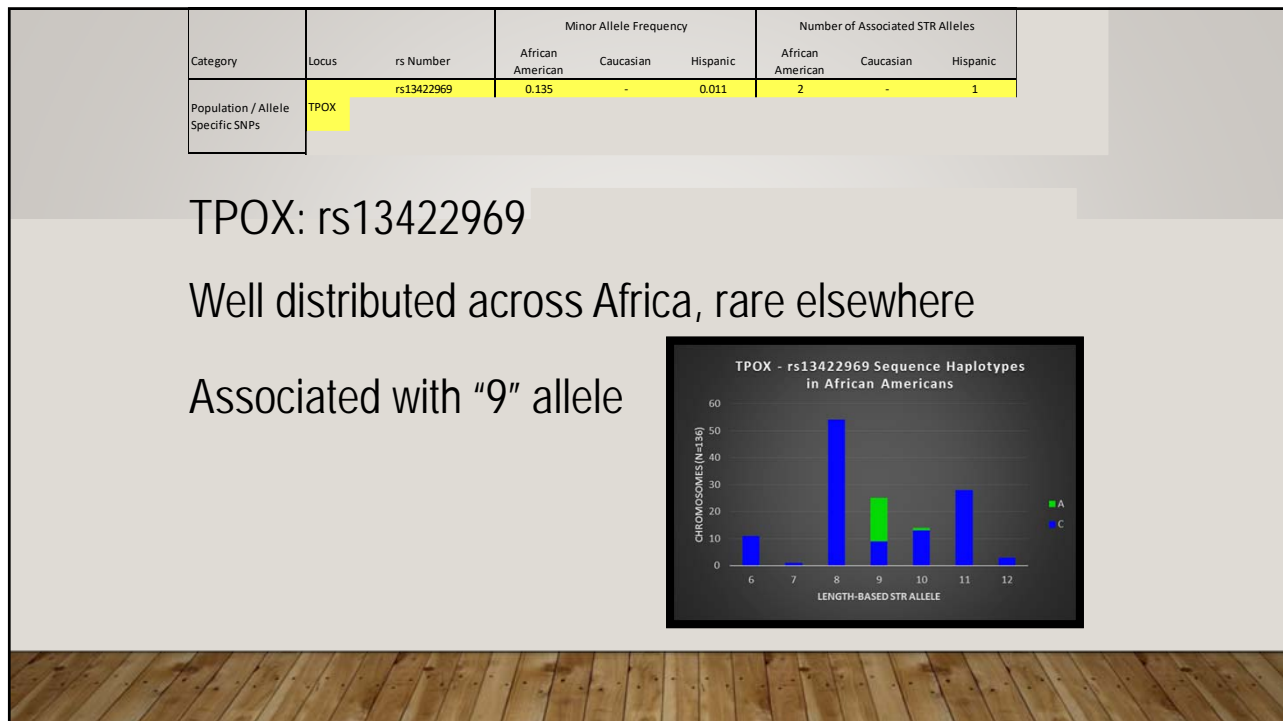
Rare SNPs were identified in the **PowerSeq** sequenced flanking regions of:

D2S441                      D18S51  
CSF1PO                     Penta D  
D8S1179                    D22S1045  
D10S1248

36



37

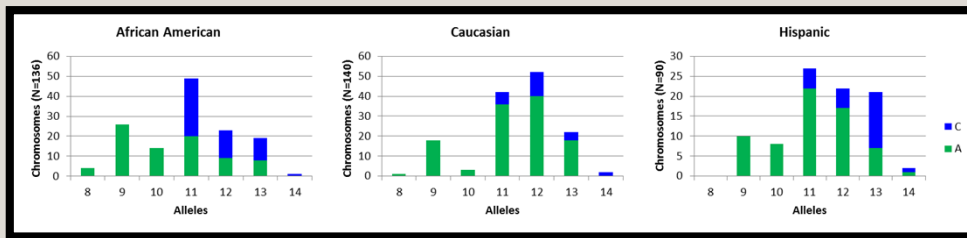


38

Category	Locus	rs Number	Minor Allele Frequency			Number of Associated STR Alleles			Alleles Gained with Flanking SNPs
			African American	Caucasian	Hispanic	African American	Caucasian	Hispanic	
Single "Old" SNP (not population or allele specific)	D16S539	rs1728369	0.404	0.171	0.278	4	4	4	5

D16S539: rs1728369

Well distributed across populations *and alleles*



39

Category	Locus	rs Number	Minor Allele Frequency			Number of Associated STR Alleles			Alleles Gained with Flanking SNPs
			African American	Caucasian	Hispanic	African American	Caucasian	Hispanic	
Multiple SNPs in Haplotype	D5S818	G-T 4bp 3'	0.316	0.257	0.189	7	5	6	17
		rs25768	0.228	0.257	0.156	5	5	5	
		rs146841551	0.029	-	-	1	-	-	
		rs541272009	0.007	-	-	1	-	-	

D5S818:  
rs73801920  
rs25768

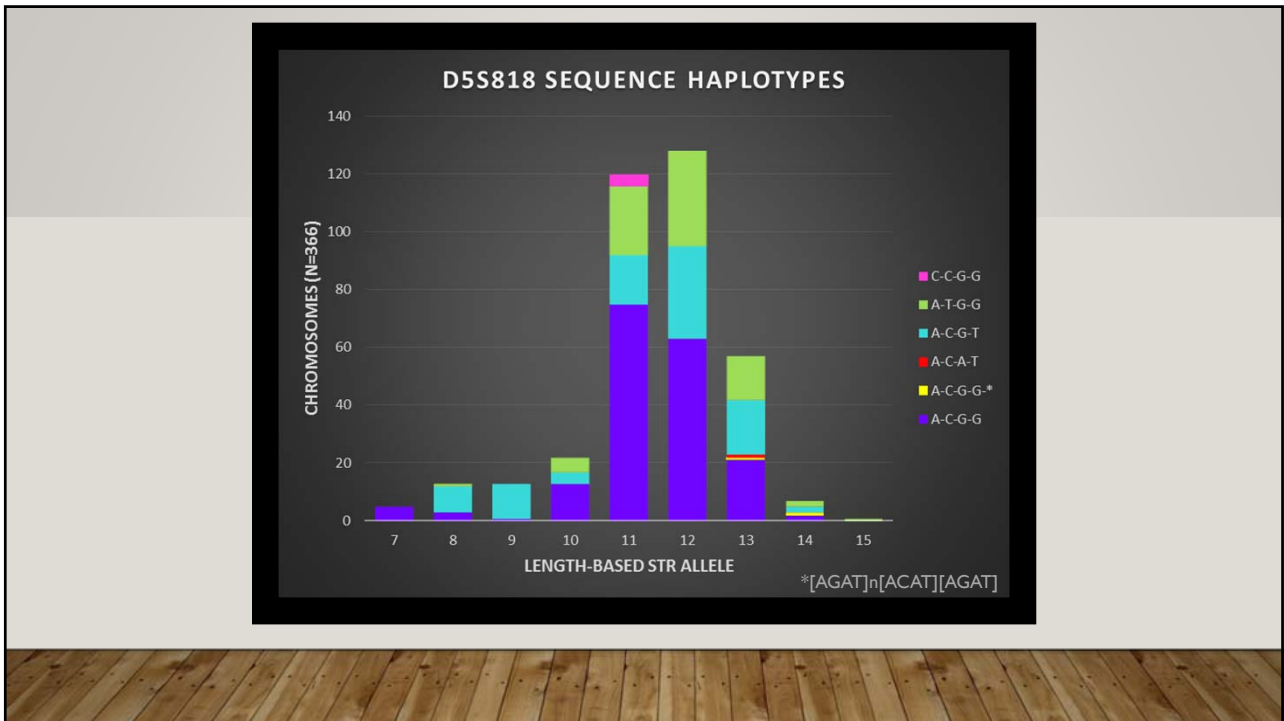
} Well distributed across populations and alleles

two additional rare SNPs

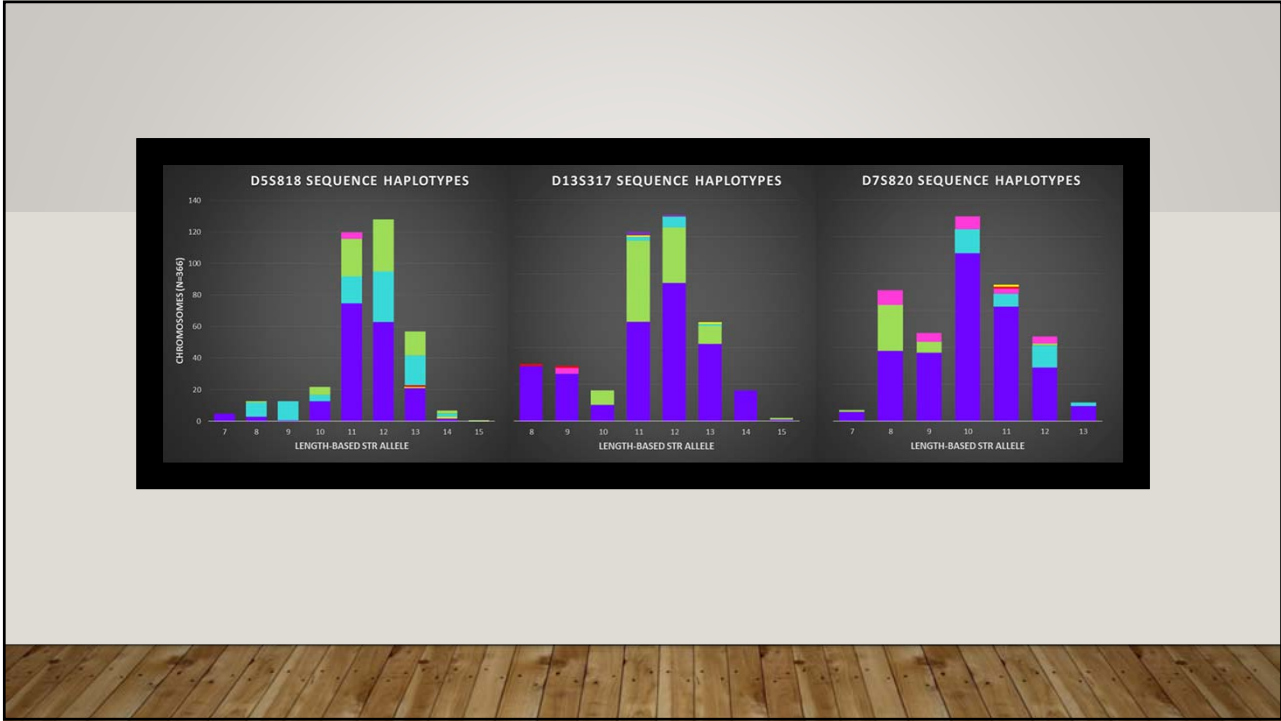
40



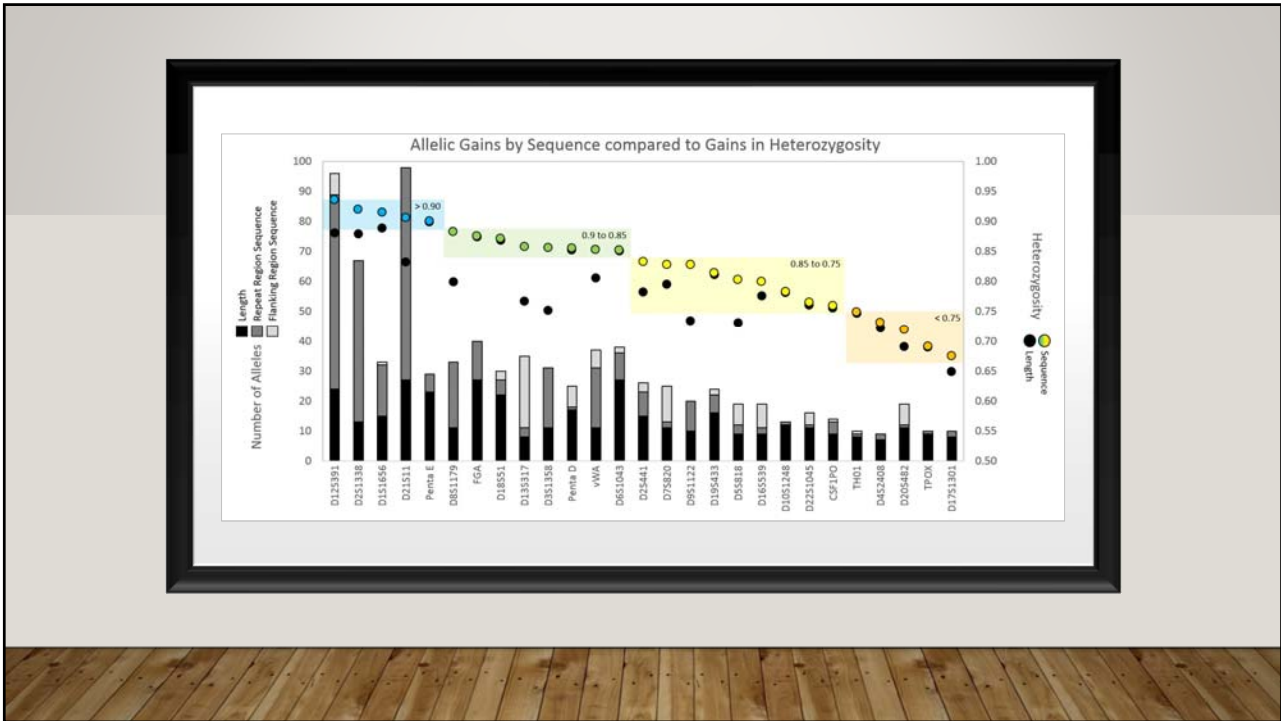
41



42

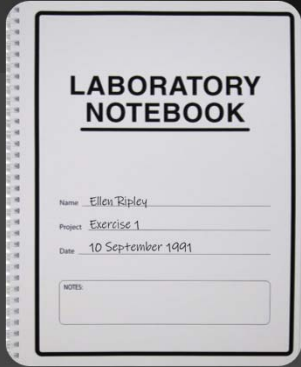


43



44

# EXERCISE I



LABORATORY  
NOTEBOOK

Name: Ellen Ripley  
Project: Exercise 1  
Date: 10 September 1991

NOTES

45

---

# HISTORICAL NOMENCLATURE

---

46

# 1991 Early paper with nomenclature info

Locus Name: based on GenBank locus designations

- e.g. HUMHPRTB

Motif: lowest alphabetical representation of the STR

- e.g. GATA > AGAT
- TH01!!!

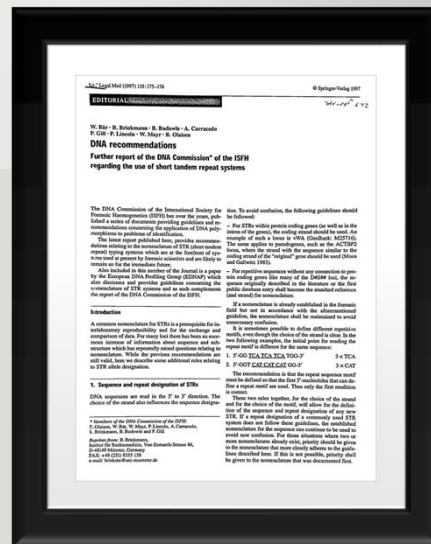
Number of repeats based on sequencing data from at least two alleles



# 1997 DNA recommendations of the ISFH

A common nomenclature for STRs is a prerequisite for interlaboratory reproducibility and for the exchange and comparison of data.

For many loci there has been an enormous increase of information about sequence and substructure which has repeatedly raised questions relating to nomenclature.





# 1997 DNA recommendations of the ISFH

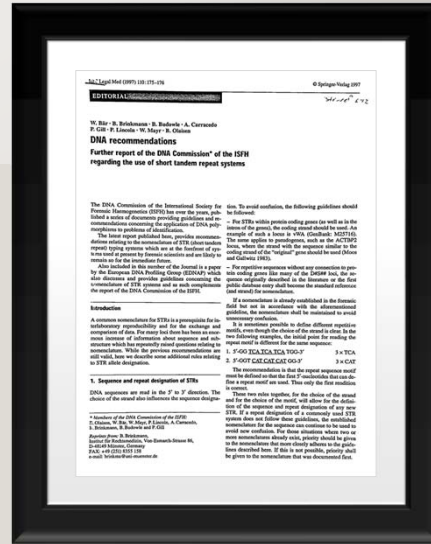
## Strand

- protein coding genes or pseudogenes: use coding strand
- no connection to coding genes: sequence originally described in literature or 1<sup>st</sup> public DB entry

## Motif (unless already established):

- First 5' nts that can define a motif

If multiple nomenclatures exist, prioritize the one which most closely adheres. If not possible, then use the first documented.



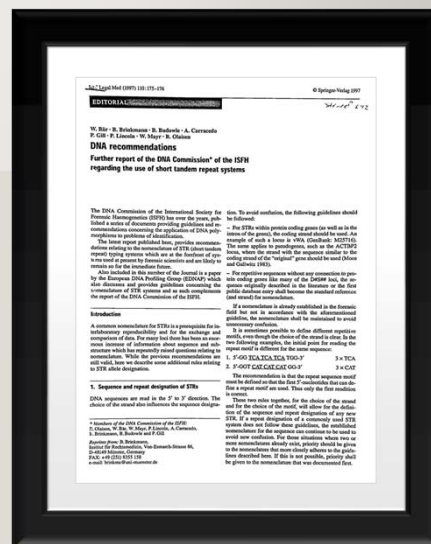
# 1997 DNA recommendations of the ISFH

## Allele Numbering for non-straightforward loci

- If sequence variation is present, use full number of repeats including variation
- Complex loci: repeat nomenclature should have a mathematical relationship to length of a consensus allele
- Example D21S11

repeats	length	relation
27	213bp	$27 \times 4 + 105^* = 213$
31	229bp	$31 \times 4 + 105^* = 229$
33.2	239bp	$33 \times 4 + 2 + 105^* = 239$

\*105bp = sum of 5'+3' flank and 43bp interspersed element



“Authors should follow these recommendations and state they did so.”

## 2006 Recommendations of the ISFG

Allele number derived from total number of contiguous variant and non-variant repeats

Single repeat units adjacent to main array and of the same sequence as main variable repeat should be included

- e.g.  $[GATA]_n [GACA]_2 [GATA] = n+2+1$

If repetitive motifs are not adjacent and have  $\leq 3$  units and show no variation, they should not be included

- e.g.  $[GATA]_n [GACA]_2 \mathbf{N}_8 [GATA]_3 = n+2$ 
  - If  $N \leq 4$  nt, include
  - If  $N > 4$  nt, exclude



51

## 2006 Recommendations of the ISFG

Indels within repeat are counted

e.g.  $[CTTTT]_8 C [CTTTT]_3 = 11.1$

Indels in flank are designated separately

e.g. 11 (U40Tins)

(this only works when sequencing)

Ambiguous indels assigned to highest numbered end of homopolymer

e.g.  $C-AAAAAAA[GATA]_6 \dots = U9A\text{del}$

e.g.  $\dots [TATC]_6 \text{TTTTTTT-GC} = D9T\text{del}$



52

## 2006 Recommendations of the ISFG

Indels within repeat are counted

- e.g.  $[CTTTT]_8 C [CTTTT]_3 = 11.1$

Indels in flank are designated separately

- e.g. 11 (U40Tins)

(this only works when sequencing)

Ambiguous indels assigned to highest numbered end of homopolymer

- e.g. C-AAAAAAA[GATA]<sub>6</sub>... = U9Adel
- e.g. ...[TATC]<sub>6</sub>TTTTTTT-GC = D9Tdel



53

## 2016 Considerations of the ISFG

Consideration I: MPS analysis should be performed with software that allows STR sequences to be exported and stored in databases as **sequence (text) strings** to capture the maximum consensus sequence information.



54

## 2016 Considerations of the ISFG

Consideration 2. The **forward strand** direction assigned in the human genome has been constant for all assemblies published since the first draft in 2001 and can be used to align STR sequences.



55

## 2016 Considerations of the ISFG

Consideration 3. The choice of reference sequence is crucial for standardizing STR nomenclature systems. At the time of writing, **GRCh38** is the most up-to-date sequence assembly and is **recommended as the framework with which to define repeat region structure** for sequence alignment and for the mapping of sequence features such as SNPs.



56

## 2016 Considerations of the ISFG

Consideration 4. Further work is needed to translate the nomenclature of STR loci thus far coded relative to the reverse strand and repeat region start and end points. There is a need to **strictly define** these and other **anchor points to specify the repeat regions.**



57

## 2016 Considerations of the ISFG

Consideration 5. Although simple STR nomenclature systems may be required at some point in the future ... comprehensive STR nomenclature systems are preferred for early adopters... **Backward compatibility to the repeat-based nomenclature derived from CE needs to be maintained** to preserve the universal applicability of established national STR databases.



58

## 2016 Considerations of the ISFG

Consideration 6. To account for relevant genetic variation outside common repeat regions, STR sequences stored as **sequence strings should include flanking sequences as well as the genome coordinates** of the sequence read start and end points.



59

## 2016 Considerations of the ISFG

Consideration 7. **Updated allele frequency databases will be necessary** to take full advantage of the increased power of discrimination offered by MPS generated STR data. A unified nomenclature system is needed to **ensure compatibility of worldwide population databases.**



60

## 2016 Considerations of the ISFG

Consideration 8. Future forensic MPS multiplexes would benefit from **retention of past markers for backward compatibility** and a marker selection process based on population data, molecular biology, sequencing chemistry, and a continued dialogue between the forensic community and commercial suppliers.



61

Fun with DI2S39I

62



Primarily:

[AGAT]<sub>6-18</sub> [AGAC]<sub>4-11</sub> [AGAT]<sub>0-1</sub>

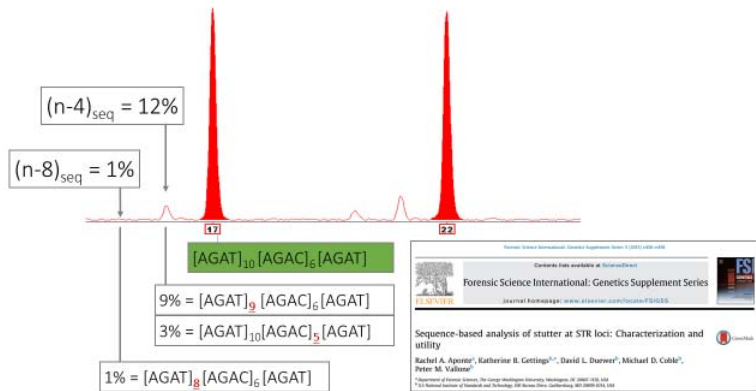
~5% Europeans:

AGAT GAT [AGAT]<sub>8-10</sub> [AGAC]<sub>7</sub> AGAT



63

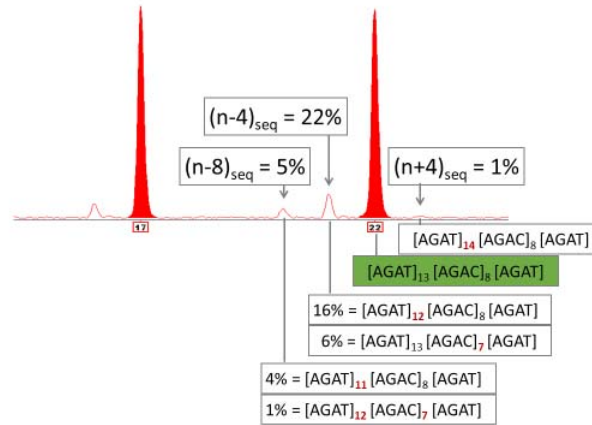
### D12S391 stutter by sequence



64



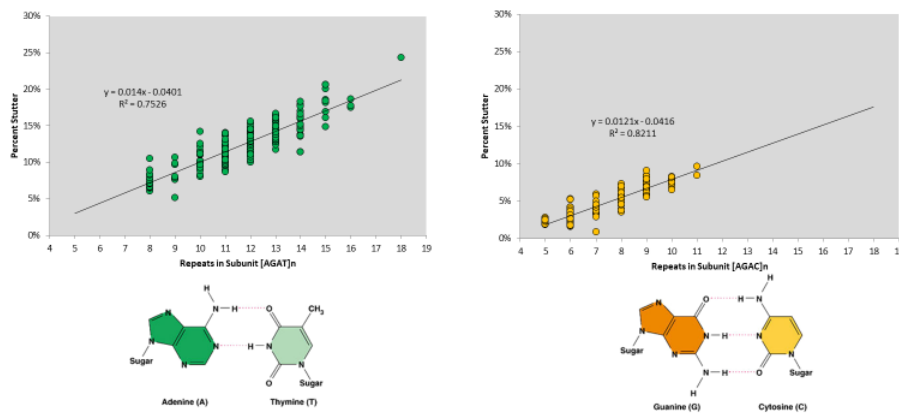
## D12S391 stutter by sequence




65

## D12S391 stutter by sequence

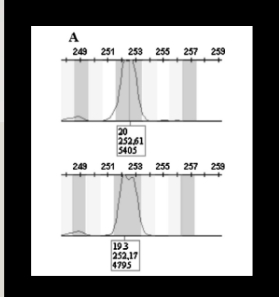
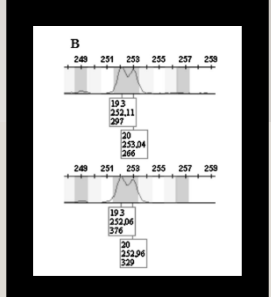
[AGAT]<sub>n</sub> trends ~3% higher stutter than [AGAC]<sub>n</sub>



66



Single Source samples show varying levels of 1bp resolution:

This has implications for validations and mixture analysis:

	Allele	Repeat Region Sequence	A	C	G	T	MW	delta	if F is labeled
Sample 1	19.3	AGAT GAT [AGAT]10 [AGAC]7 AGAT	39	7	20	13	10709.93		
Sample 2	20	[AGAT]13 [AGAC]6 AGAT	40	6	20	14	10860.08	150.15	easier
Sample 3	20	[AGAT]12 [AGAC]8	40	8	20	12	10830.04	120.11	harder

67

# STRBase & STRSeq

68

# STRBase

## Short Tandem Repeat DNA Internet DataBase

69

The collage consists of three overlapping documents:

- Left Document:** A research paper titled "STRBase: a short tandem repeat DNA database for the human identity testing community" by Christian M. Ruitberg, Dennis J. Reeder, and John M. Butler. It includes an abstract, introduction, and contact information.
- Middle Document:** The main STRBase website page, featuring the title "STRBase Short Tandem Repeat DNA Internet DataBase" and sections for "General Information" and "Forensic Interest Data".
- Right Document:** A technical report for the D8S1179 locus, showing "General Information", "PCR Primer Information", and a table of "PCR Product Sizes of Observed Alleles".

70

The screenshot displays the STRBase website interface. At the top, the NIST logo and the text "STRBase: A web based resource for STR information" are visible. The page is divided into several columns. On the left, there are navigation links for "Commonly Used Auto STRs", "Other Auto STRs", "X-Chromosome STRs", and "Y-Chromosome STRs". The main content area includes sections for "Introduction", "Purpose", "New Features", "Updates", "Navigation", "Search", and "Tables". A search bar is located at the top right. Below the main content, there are several charts and graphs, including a pie chart and a bar chart. The bottom of the page features the NIST Applied Genomics logo and contact information for the Applied Genomics Group.

71

The graphic consists of a light blue rectangular area with a white border. In the center, the word "STRSeq" is written in a large, bold, black sans-serif font. Below it, the text "The STR Sequencing Project" is written in a smaller, black sans-serif font. The bottom of the graphic features a horizontal strip with a wooden floor texture.

72



73

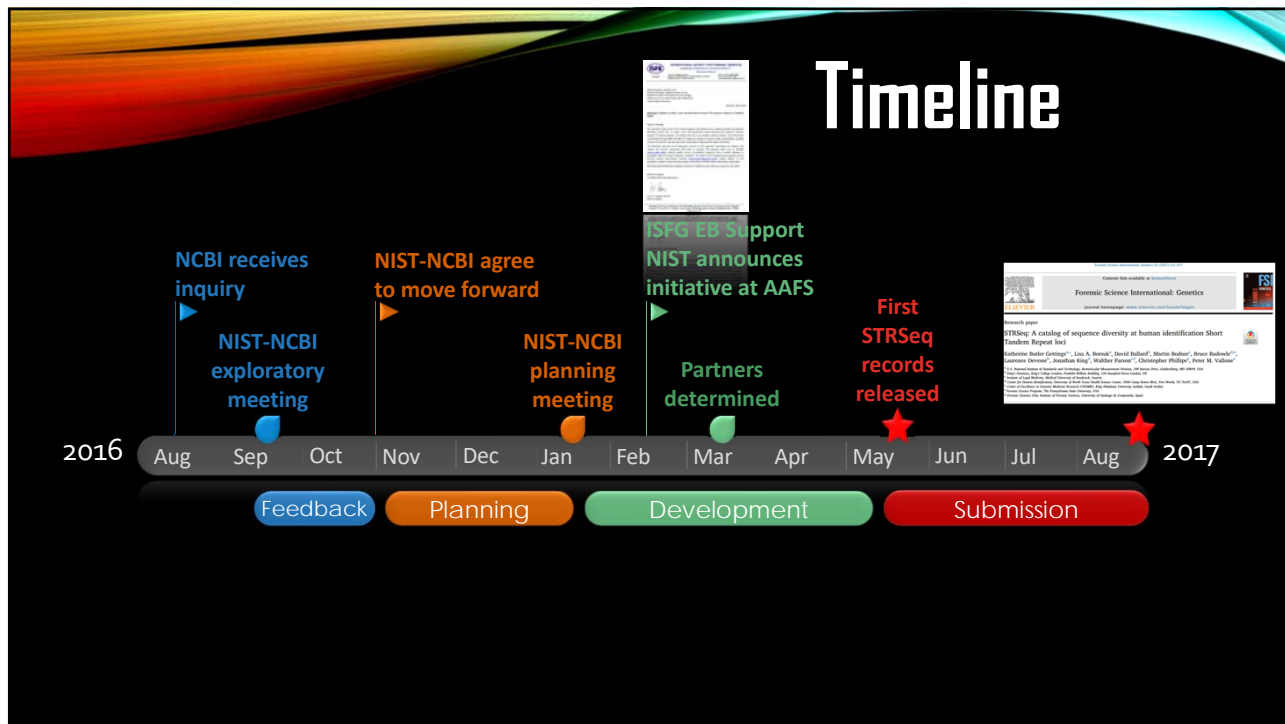
A title slide for a presentation. The background features a dark blue and purple space theme with a glowing star. A silhouette of a world map is overlaid on the background, with two large red curved arrows forming a circle around it. The word "Inception" is written in a large, white, sans-serif font at the top right. At the bottom, there are two informational boxes. The first box is titled "RefSeq: NCBI Reference Sequence Database" and contains the text: "A comprehensive, integrated, non-redundant, well-annotated set of reference sequences including genomic, transcript, and protein." The second box is titled "LocusReferenceGenomic" and contains the text: "LRG sequences provide a stable genomic DNA framework for reporting mutations with a permanent ID and core content that never changes."

**Inception**

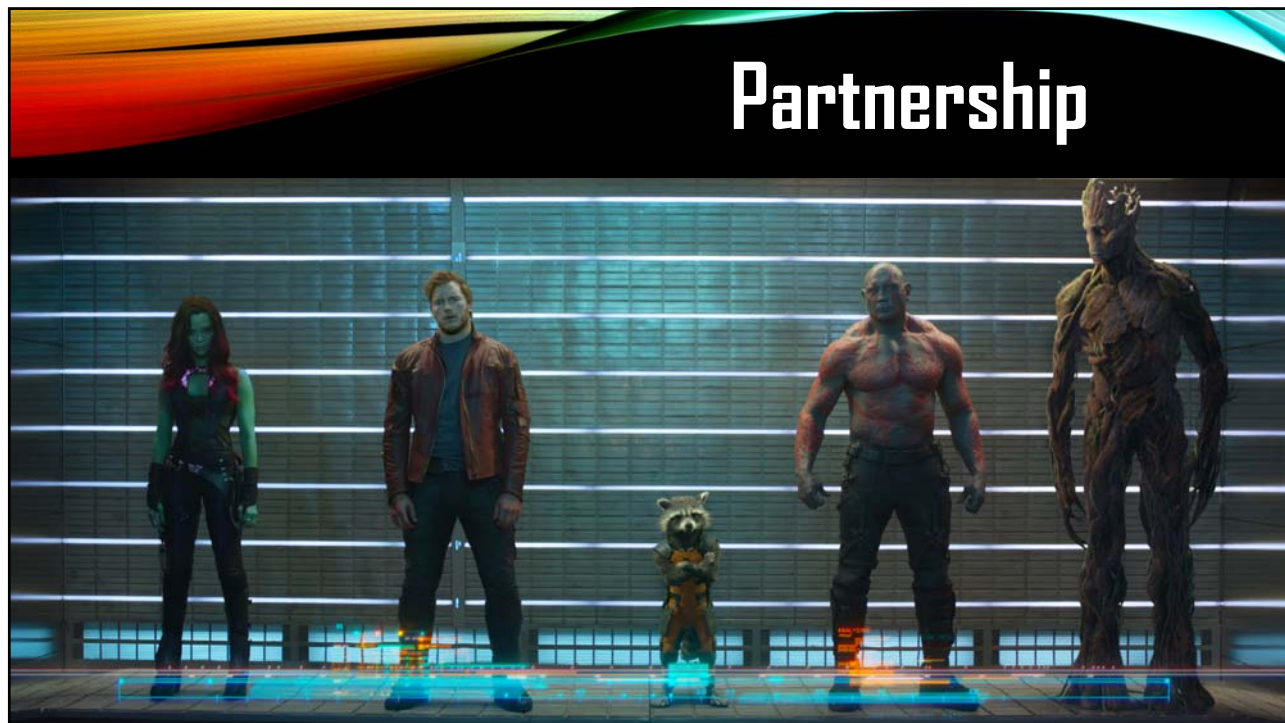
**RefSeq: NCBI Reference Sequence Database**  
A comprehensive, integrated, non-redundant, well-annotated set of reference sequences including genomic, transcript, and protein.

**LocusReferenceGenomic**  
LRG sequences provide a stable genomic DNA framework for reporting mutations with a permanent ID and core content that never changes.

74



75



76

# Partners-Roles

**UNT HEALTH SCIENCE CENTER**

**NIST**  
Population Samples  
Project Coordination  
Record Submission

**KING'S College LONDON**

**UNIVERSITÄT TIROL**  
Project Input  
STRidER Integration

**NCBI**  
Project Input  
Hosting

**USC**  
UNIVERSIDAD DE SANTIAGO DE COMPOSTELA  
Population Data  
Project Input

77

# STRSeq Samples

1786 ForenSeq  
+  
650 PowerSeq  
+  
CE supporting data

1043 ForenSeq  
+  
CE supporting data

839 ForenSeq

944 ForenSeq

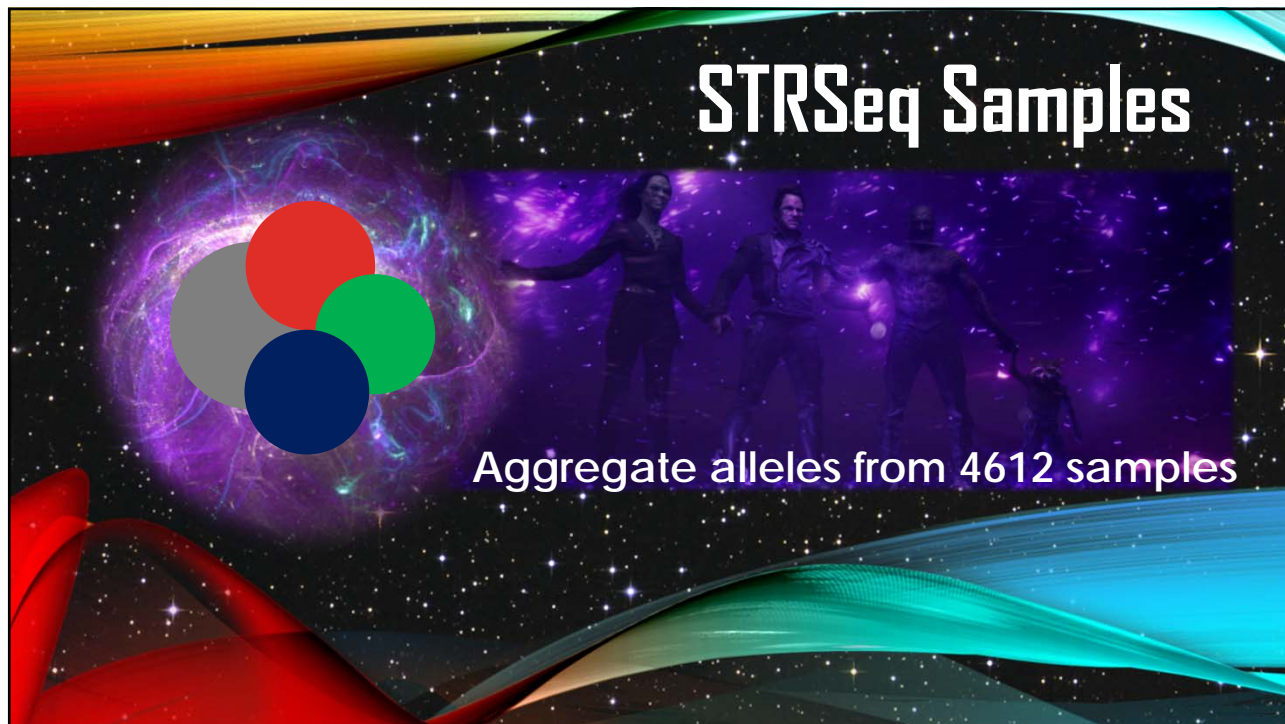
**NIST**

**KING'S College LONDON**

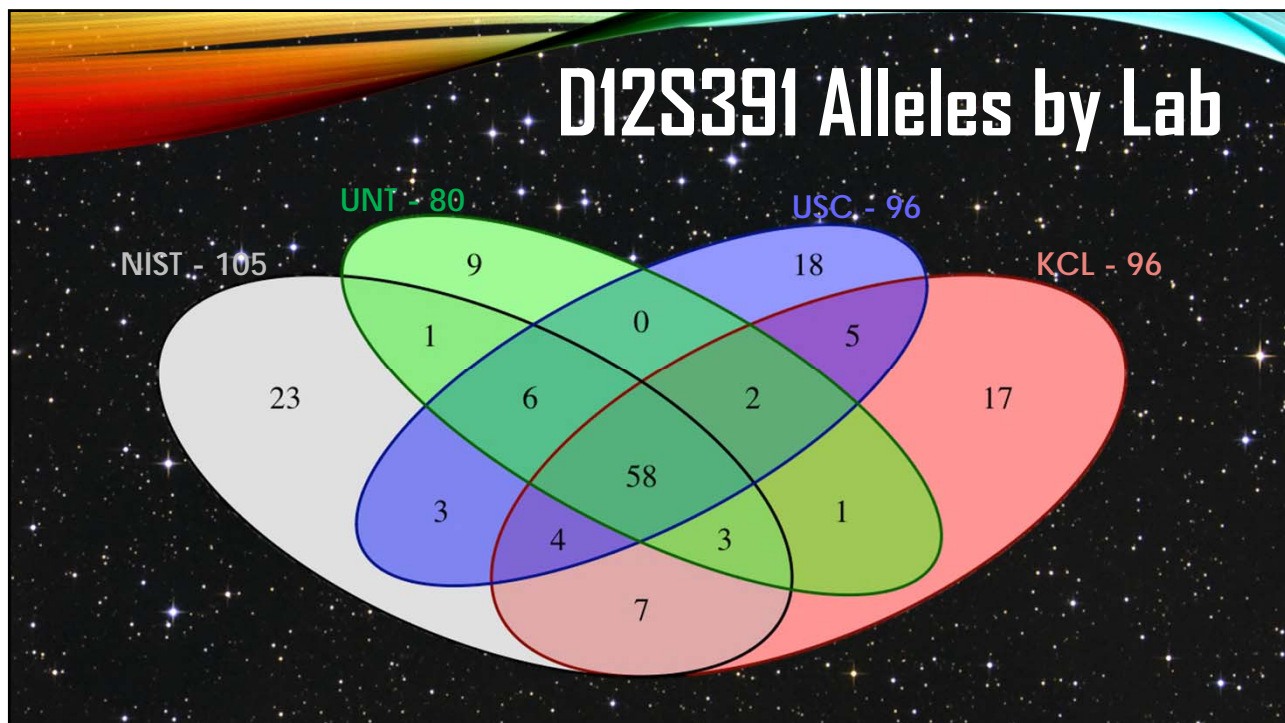
**UNT HEALTH SCIENCE CENTER**

**USC**  
UNIVERSIDAD DE SANTIAGO DE COMPOSTELA

78



79



80



https://www.ncbi.nlm.nih.gov/bioproject/380127

The screenshot displays the NCBI BioProject page for 'The STR Sequencing Project (human)'. The page includes a search bar, navigation tabs, and a main content area with the following sections:

- The STR Sequencing Project (human)**: Accession: PRJNA380127 ID: 380127. Description: The purpose of STRseq is to facilitate the description of sequence-based alleles at the Short Tandem Repeat (STR) loci targeted in human identification assays.
- Project Data**: A table showing the number of links for various resource names. Two red arrows point to the 'Nucleotide (Genomic DNA)' (1442) and 'PubMed' (5) rows.
- The STR Sequencing Project (human) encompasses the following 4 sub-projects:** A table listing sub-projects with columns for Project Type, BioProject, Name, and Title.

81

The screenshot displays the NCBI GenBank entry for 'Homo sapiens microsatellite TPOX 7 [AATG]7 rs115644759 sequence'. The page includes the following sections:


- GenBank: MF044247.1**: FASTA format.
- FASTA**: The sequence data in FASTA format.
- Repeat region**: A detailed view of the microsatellite repeat region, showing the sequence and its position within the genomic context.


82




# STRSeq in Bioinformatics

## Standalone and API BLAST

 **Download BLAST**  
Get BLAST databases and executables

 **Use BLAST API**  
Call BLAST from your application

 **Use BLAST in the cloud**  
Start an instance at a cloud provider

## Embedding the NCBI Sequence View in Web Content

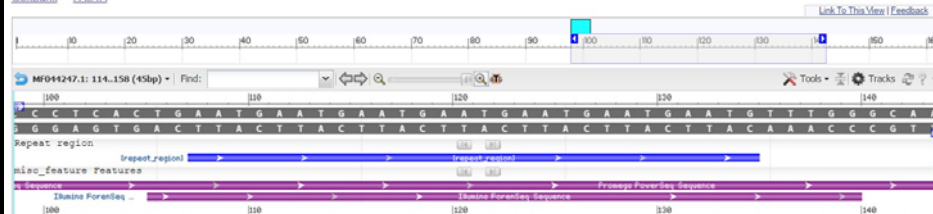
### Introduction

The NCBI Graphical Sequence Viewer (SV) is a general purpose tool for viewing biological sequence data. The Sequence Viewer has a very rich set of options and can display virtually any sequence. It can be embedded in a wide variety of web pages serving many different needs. This page has examples showing best practice for embedding Sequence Viewer with several different sets of options.

### Homo sapiens microsatellite TPOX 7 [AATG]7 rs115644759 sequence

GenBank: MF044247.1

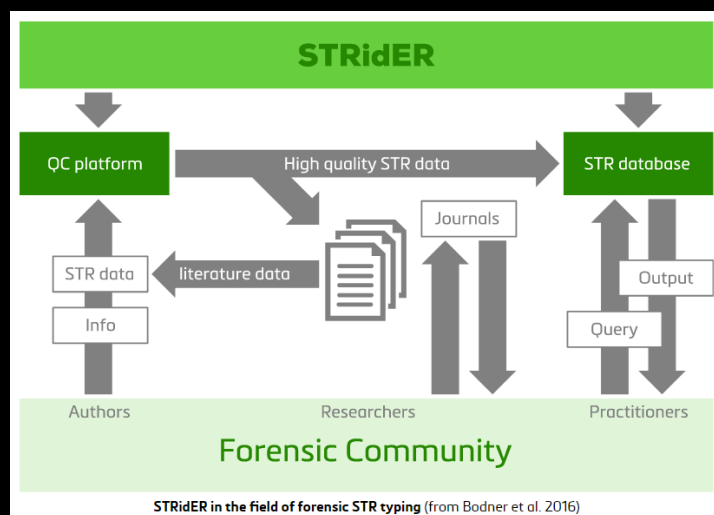
[GenBank](#) [FASTA](#)



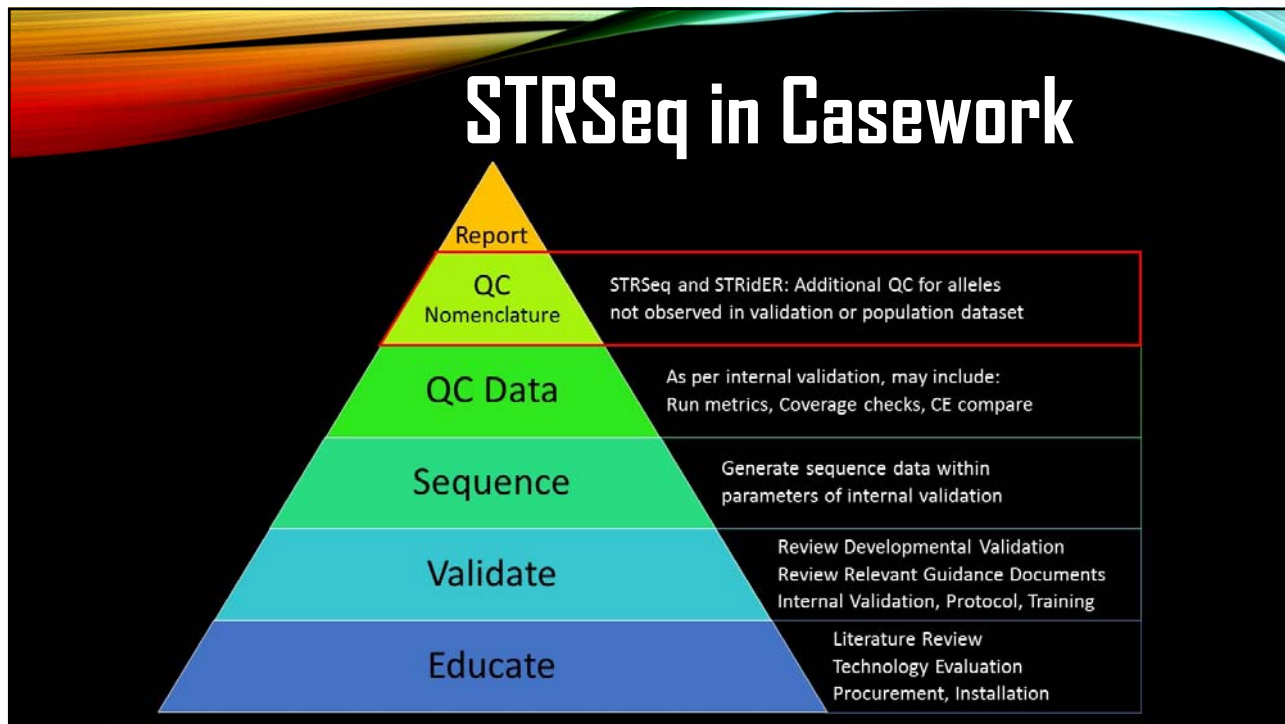
85

# STRSeq in Population Data

STRSeq



86



87

## FAQs:

How do you decide which loci to include?  
 Why do some loci not have records yet?  
 Where does the data come from?  
 Where are the allele frequencies?

Can I send you a sequence?

STRSeq			
NCBI BioProject: PRJNA380127		strseq.nist.gov	
Commonly Used Autosomal STR Loci	Alternate Autosomal STR Loci	Y-Chromosomal STR Loci	X-Chromosomal STR Loci
D1S1656	1	DYF387S1	DYS458
TPOX, D2S441, D2S1338	2	DYS19	DYS460
D3S1358	3	DYS385a-b	DYS461
FGA	4	DYS389I	DYS481
D5S818, CSF1PO	5	DYS389II	DYS505
SE33, D6S1043	6	DYS390	DYS522
D7S820	7	DYS391	DYS533
D8S1179	8	DYS392	DYS549
	9	DYS393	DYS570
D10S1248	10	DYS437	DYS576
TH01	11	DYS438	DYS612
VWA, D12S391	12	DYS439	DYS635
D13S317	13	DYS448	DYS643
	14	DYS456	Y-GATA-H4
Penta E	15		
D16S539	16		
	17		
D18S51	18		
D19S433	19		
	20		
D21S11, Penta D	21		
D22S1045	22		

88

**STRSeq: A resource for sequence-based STR analysis**

**Method:** The STR Sequencing Project (STRSeq) was initiated to facilitate the description of sequence-based alleles at Short Tandem Repeat (STR) loci required in forensic identification assays. STRSeq data are managed as GenBank records at the U.S. National Center for Biotechnology Information (NCBI). Each GenBank record contains: (1) observed sequence of an STR region, (2) annotation of the repeat region ("RepeatSeq") consistent with the guidelines of the International Society for Forensic Genetics and Forensic region nomenclature, (3) information regarding the sequencing assay and allele quality, and (4) forward complete length-based allele designation. STRSeq GenBank records are organized within a BioProject at NCBI (<https://www.ncbi.nlm.nih.gov/bioproject/26027>), which is available to:

- Chromosomal autosomal STR loci
- Alternative autosomal STR loci
- X-chromosomal STR loci
- Y-chromosomal STR loci

**Current participating laboratories and samples provided for evaluation:** NIST, KINSHIP, UNT HEALTH, USC.

**Method:** Data from three Forensic STR sequencing labs are currently being evaluated for inclusion in STRSeq: (1) Forensic DNA Signature Prep kit (Forensic), (2) PowerPlex NGP (Forensic), and (3) GlobalFiler HD (Forensic). All of the samples currently represented in STRSeq have been sequenced with ForensicSeq, NIST's sequencer-based STR analysis pipeline. The current STRSeq records include all of the STR loci and alleles currently being sequenced with GlobalFiler HD, which are currently being reviewed for inclusion in the STRSeq and GlobalFiler HD across the sequencing sites.

**What am I looking at?**  
 STRSeq resources defined in a hierarchical representation of the STR Sequencing Project (STRSeq), which is divided into four subgroups: (1) Commonly used autosomal STR loci, (2) Alternate autosomal STR loci, (3) X-chromosomal STR loci, and (4) Y-chromosomal STR loci. Currently, only the autosomal STR subgroups contain sequence records. Colors within the loci represent observed length-based alleles, which contain colored allele designations.

**What do the colors represent?**  
 Colors designate the laboratory from which the unique sequence was identified:  
 Multiple Labs (grey), NIST (red), KIN (blue), UNT (green), USC (purple).

89

EXERCISE 2

LABORATORY NOTEBOOK

Name: Ellen Ripley


Project: Exercise 2

Date: 10 September 2019

NOTES:

90

**COFFEE BREAK**  
meet back in 30 minutes

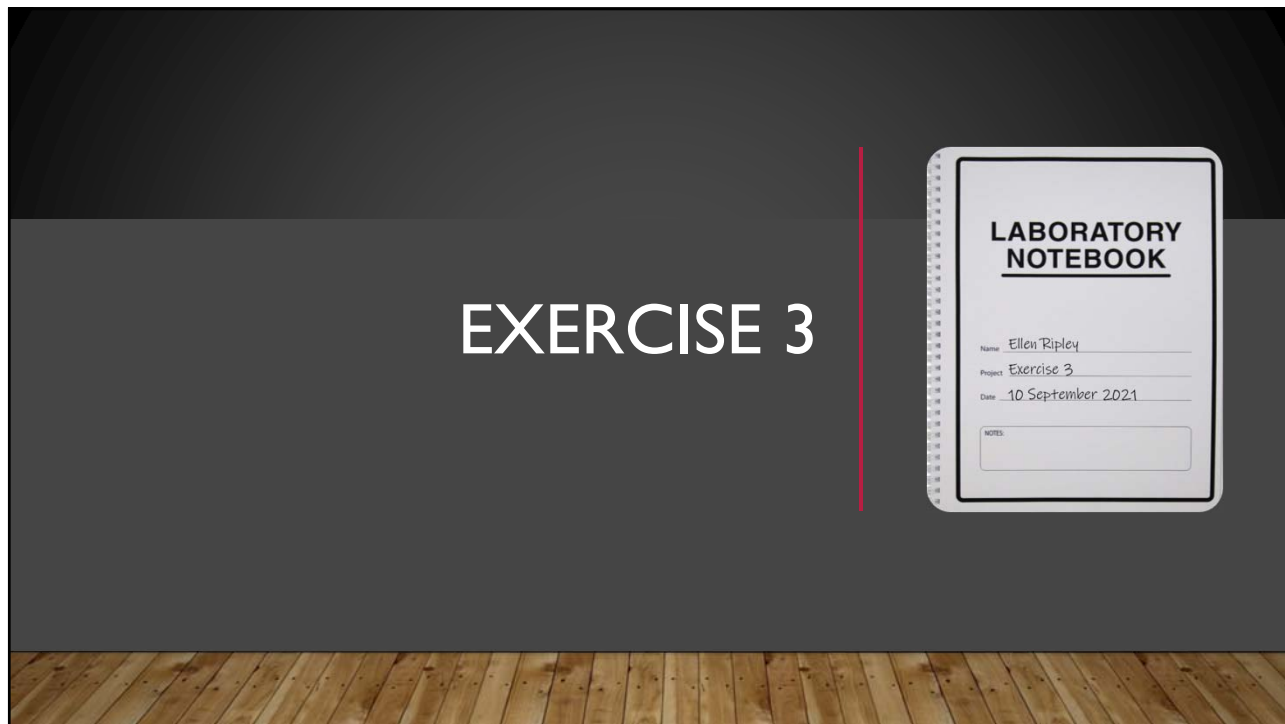
A close-up photograph of a coffee cup with a smiling face drawn on the foam. The face has two small dots for eyes and a curved line for a mouth. The coffee is dark, and the foam is light brown. The cup is on a wooden surface.

91

**NEW IDEAS IN  
NOMENCLATURE**

**SID & STRNaming**

92



93



94



# STRAND *working group*

align | name | define

95

# STRAND *working group*

align | name | define



Our mission is to harmonize related efforts across member laboratories:



▶ STRidER STR sequence quality control



▶ STRSeq catalog of sequences



▶ STRaitRazor bioinformatic freeware



▶ Forensic STR Sequence Guide

and to characterize additional STR loci present in the genome which may be useful for forensic purposes in the future.


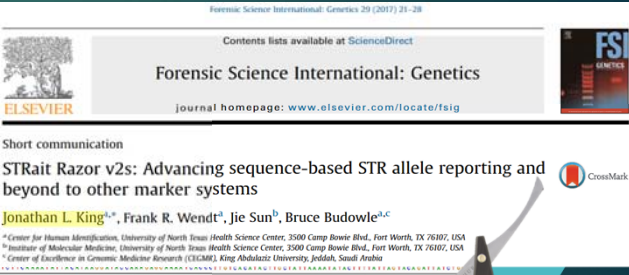
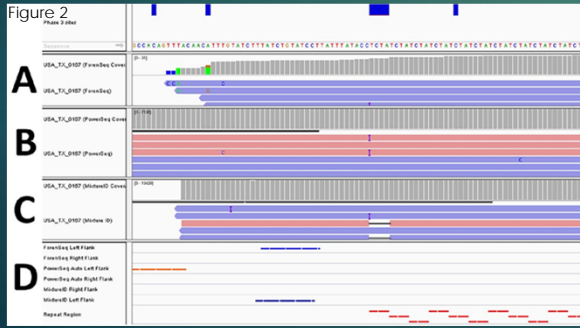
96



# STRAND *working group*

align | name | define

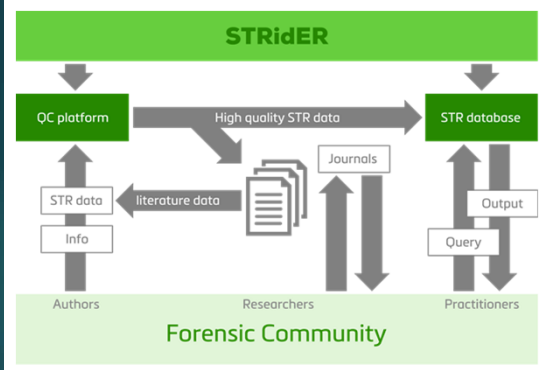
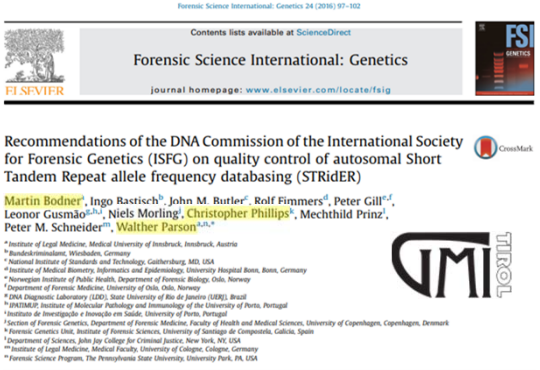
► STRait Razor Agnostic Freeware

97

# STRAND *working group*

align | name | define

The existing architecture of STRidER allows for the implementation of nucleotide sequence strings and thus is fully compatible with the QC of population data generated by MPS.

strider.online

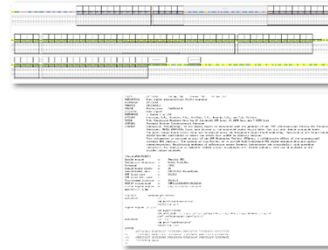
98

# STRAND *working group*

align | name | define

## QC of STR sequence data on STRidER

- submission of FASTA-like strings per locus per individual
- alignment and comparison to references:
  - Forensic STR Sequence Structure Guide**
  - STRSeq catalogue**
- translation into CE allele
- inspection of flanking region where available
- **no STR sequence nomenclature** required, assessed or returned in reports
- database will host STR sequence data in future

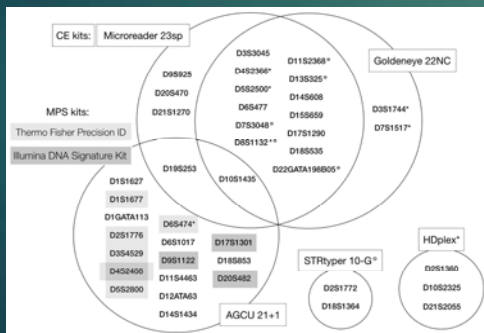


99

# STRAND *working group*

align | name | define

## Forensic STR Sequence Guide

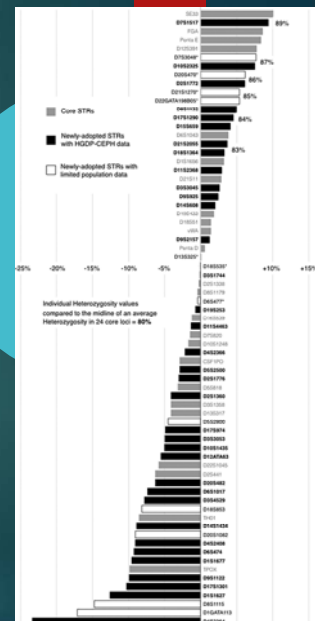


Forensic Science International: Genetics

Research paper: A genomic audit of newly-adopted autosomal STRs for forensic identification

C. Phillips

Abstract: An audit of autosomal STR multiplex kits has recently been completed for the commonly used 24 'core' forensic STRs... The results of this audit are presented in this paper...




100

# STRAND *working group*

align | name | define

► Forensic STR Sequence Guide



Forensic Science International: Genetics 34 (2018) 162–169

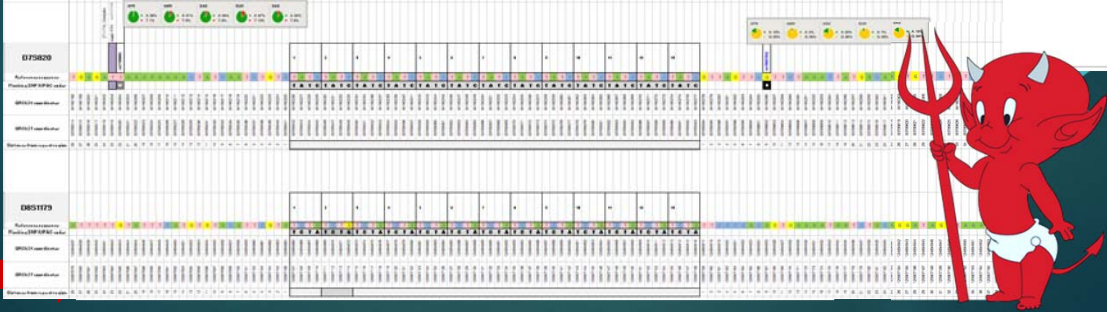
Contents lists available at ScienceDirect

Forensic Science International: Genetics

journal homepage: [www.elsevier.com/locate/fsigen](http://www.elsevier.com/locate/fsigen)

“The devil’s in the detail”: Release of an expanded, enhanced and dynamically revised forensic STR Sequence Guide

C. Phillips<sup>a,\*</sup>, K. Butler Gettings<sup>b</sup>, J.L. King<sup>c</sup>, D. Ballard<sup>d</sup>, M. Bodner<sup>e</sup>, L. Borsuk<sup>b</sup>, W. Parson<sup>d,f</sup>



101

## STR Nomenclature Meeting

April 11-12, 2019 London



5' to 3':  
 Walther Parson, Lisa Borsuk,  
 Peter Schneider, Brian Young,  
 Rebecca Just, Jodi Irwin, David  
 Ballard, Sascha Willuweit, Cydne  
 Holt, Chris Phillips, Jonathan King,  
 Tunde Huszar, Peter Gill, Christian  
 Sell, Kris Van der Gaag, Laurence  
 Devesse, Claus Borsting, Doug  
 Hares, Katherine Gettings, Rob  
 Lagace, Jerry Hoogenboom,  
 Martin Bodner, Peter deKnijff,  
 Sebastian Ganschow, Pedro  
 Barrio, Teresa Gross

STRAND *working group*

D.Ballard:KCL,UK | M.Bodner&W.Parson:GVL,Austria | L.Borsuk&K.Gettings:NST,US | J.King:UNT,US | C.Phillips:USC,Spain

102

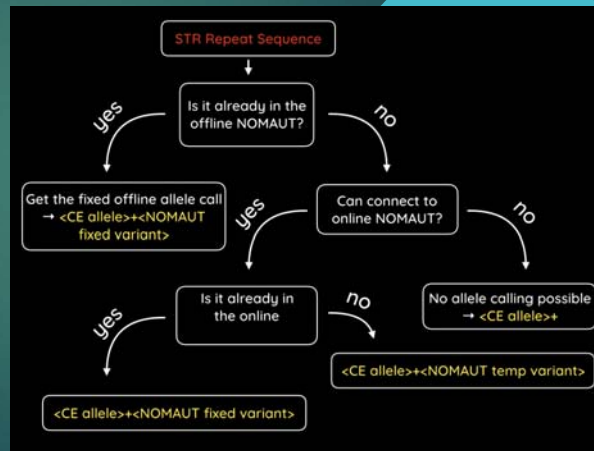
# STR Nomenclature Meeting

## Formats for STR Sequences

▶ Short Designator



Sascha Willuweit  
Charité Medical University

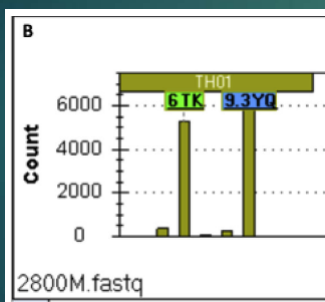


103

# STR Nomenclature Meeting

## Formats for STR Sequences

Short Designator



Adapted from Figure 1

Forensic Science International: Genetics 42 (2019) 14–20

Contents lists available at ScienceDirect

**Forensic Science International: Genetics**

Journal homepage: [www.elsevier.com/locate/bsign](http://www.elsevier.com/locate/bsign)

Research paper

**A nomenclature for sequence-based forensic DNA analysis**

Brian Young<sup>a</sup>, Tom Faris, Luigi Armogida

NicheVision Forensics, LLC, 526 South Main St. Akron, OH, 44311, USA

Adapted from Table 2

Sequence Type	Read Count	Sequence	SID Code
Allele	5,272	TCCATGGTGAATGAATGAATGAATGAATGAATGAGGGAATAAGG	6 TK
Sequence Artifact	35	TCCATGGTGAATGAAGGAATGAATGAATGAATGAGGGAATAAGG	6 VS
N-1 Stutter	368	TCCATGGTGAATGAATGAATGAATGAATGAGGGAATAAGG	5 14
Allele	6,653	TCCATGGTGAATGAATGAATGAATGAATGAATGAATGAATGAGGGAATAAGG	9.3 YQ
Sequence Artifact	11	TaCAtGGTGAATGAATGAATGAATGAATGAATGAATGAGGGAATAAGG	9.3 MS
Sequence Artifact	11	TCCATGGTGAATGAATGAATGAATGAATGAATGAATGAGGGAATAAGG	9.3 ZK
N-1 Stutter	232	TCCATGGTGAATGAATGAATGAATGAATGAATGAATGAGGGAATAAGG	8.3 WC

104

# STR Nomenclature Meeting

## Formats for STR Sequences

### Short Designator

#### Longest Uninterrupted Stretch (LUS) Concept for Sequence Allele Designation

##### Was not proposed as sequence allele nomenclature

- Aim was representation of sequence alleles for probabilistic genotyping
- Practical, near-term solution for labs already employing prob gen for casework
  - Lack of program options for sequence alleles could be a barrier to NGS adoption
  - Could be used in mixture *interpretation* even prior to storage of sequence alleles in databases



#### Example D8S1179 alleles

Length (RU)	Bracketed format	LUS length	LUS Allele Designation (RU_LUS)
12	[TCTA] <b>12</b>	<b>12</b>	<b>12_12</b>
12	[TCTA]2 TCTG [TCTA] <b>9</b>	<b>9</b>	<b>12_9</b>
12	TCTA TCTG [TCTA] <b>10</b>	<b>10</b>	<b>12_10</b>

UNCLASSIFIED//FOUO

3

105

# STR Nomenclature Meeting

## Formats for STR Sequences

### Bracketed Repeat

#### STRNaming from NFI

- ▶ Jerry Hoogenboom & Kris van der Gaag



```

CE11_TATC[8]TGTC[1]TATC[3]AATC[1]ATCT[3]
CE11_TATC[10]AATC[3]ATCT[3]
CE11_TATC[11]AATC[2]ATCT[3]
CE11_TATC[12]AATC[1]ATCT[3]
CE11_TATC[12]AATC[1]ATCT[3]_-24G>A
CE11_TATC[12]AATC[1]ATCT[3]_-25C>T
CE11_TATC[13]ATCT[3]
CE12_TATC[7]TATT[1]TATC[5]AATC[1]ATCT[3]
CE12_TATC[12]AATC[2]ATCT[3]
CE12_TATC[13]AATC[1]ATCT[3]
CE12_TATC[13]AATC[1]ATCT[3]_-24G>A
CE12_TATC[13]AATC[1]ATCT[3]_-25C>T
CE12_TATC[13]AATC[2]ATCT[2]
CE12_TATC[14]ATCT[3]
CE13_TATC[13]AATC[2]ATCT[3]
CE13_TATC[14]AATC[1]ATCT[3]
CE13_TATC[14]AATC[1]ATCT[3]_-24G>A
CE13_TATC[14]AATC[1]ATCT[3]_-25C>T
CE13_TATC[15]AATC[1]ATCT[3]_+9GTCT>
CE13_TATC[15]ATCT[3]
  
```

106

## Flanking Region Polymorphism Nomenclature

Option	Description	Example	Requires Range	Requires Reference Genome	Requires InDel Alignment Parameters
1	Report differences relative to a reference genome	Chr5:123775552:C>A	Yes	Yes	Yes
2	Report differences relative to repeat region	+4C>A	Yes	Sort of	Yes
3	Report rs number and change	rs73801920 C>A	Yes	No, but	No

107

## STR Nomenclature Meeting

### Formats for STR Sequences

#### Full String = Unequivocal Record

Storage method/location is lab-determined

*"At this time, forensic DNA databasing software (e.g. CODIS) is generally not equipped to store or search STR sequence strings.*

*Such databases primarily contain convicted offender samples; therefore, enabling STR sequence storage or search capabilities may be of limited use until laboratories begin routinely sequencing this sample type.*

*In the interim, length based (numerical allele) profiles can be developed via STR sequencing assays.*

*Analysts confirming interlaboratory matches could compare sequence data, when applicable."*

108

# STR Nomenclature Meeting

## Defined Coordinates

Assay Specific

109

# STR Nomenclature Meeting

## Defined Coordinates

PowerSeq 46GY	GeneMarker NGS Range
ForenSeq DNA Signature Prep Kit	UAS Flanking Region Report Range
Precision ID GlobalFiler NGS v2	Converge .bed file range

Supplementary File - 24 auSTRs


110

# STR Nomenclature Meeting

## Forensic Specific Reference


- ▶ Advantages
  - ▶ Elimination of rare SNP alleles in STR flanking regions, incorporation of known insertions
  - ▶ Stability - the forensic
- ▶ Disadvantages
  - ▶ Significant effort would be required for curation, maintenance, version control, and enforcement of general use within the

**AFR**



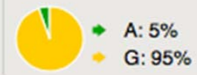
A: 18%  
G: 82%

**AMR**




A: 17%  
G: 83%

**ASN**



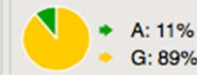
A: 5%  
G: 95%

**EUR**



A: 27%  
G: 73%

**SAS**



A: 11%  
G: 89%


of worldwide populations, or representative of maximal complexity

bioinformatic methods

111

# STRAND *working group*

align | name | define



Forensic Science International: Genetics 22 (2016) 54-63

Contents lists available at ScienceDirect

ELSEVIER

journal homepage: [www.elsevier.com/locate/fgs](http://www.elsevier.com/locate/fgs)

Massively parallel sequencing of forensic STRs: Considerations of the DNA commission of the International Society for Forensic Genetics (ISFG) on minimal nomenclature requirements

Walther Parson<sup>a,b,\*</sup>, David Ballard<sup>c</sup>, Bruce Budowle<sup>d,e</sup>, John M. Butler<sup>f</sup>, Katherine B. Gettings<sup>g</sup>, Peter Gill<sup>h,i</sup>, Leonor Gusmão<sup>j,k</sup>, Douglas R. Hares<sup>l</sup>, Jodi A. Irwin<sup>l</sup>, Jonathan L. King<sup>g</sup>, Peter de Knijff<sup>m</sup>, Niels Morling<sup>n</sup>, Mechthild Prinz<sup>o</sup>, Peter M. Schneider<sup>o</sup>, Christophe Van Neste<sup>o</sup>, Sascha Willuweit<sup>o</sup>, Christopher Phillips<sup>o</sup>

Our mission is to harmonize related efforts across member laboratories:

- ▶ STRidER STR sequence quality control
- ▶ STRSeq catalog of sequences
- ▶ STRaitRazor bioinformatic freeware
- ▶ Forensic STR Sequence Guide

and to characterize additional STR loci present in the genome which may be useful for forensic purposes in the future.

112



# STRAND *working group*

align | name | define



## Genome in a Bottle (GIAB)

### 7 Coriell cell lines

- ▶ One individual and two trios
- ▶ PCR-free prep, HiSeq, PacBio, ONT, 10X

### Analyzing STR regions in GIAB samples

- ▶ Any "novel" marker can be characterized
- ▶ Proof of concept targeting >600 STRs



**Unleashing Novel STRs**  
via characterization of  
**Genome in a Bottle**  
reference samples

**What is GIAB?**  
A resource named by I2D established to **AUTHORITATIVE CHARACTERIZATION** of benchmark human genomes.

**How do we extract STR data from GIAB genomes?**  
Align with target → Extract data from GIAB using target Coordinates → Custom analysis for flanking and flanking data files → Accept results of multiple analyses

**How can we use this resource?**  
STR Marker Discovery/Evaluation → QC for novel STR assay targets → Input for nomenclature discussions

113

# QUESTIONS?

enjoy your week at the  
28<sup>th</sup> ISFG in Prague

Questions later? [katherine.gettings@nist.gov](mailto:katherine.gettings@nist.gov) [strseq@nist.gov](mailto:strseq@nist.gov) [strandwg@gmail.com](mailto:strandwg@gmail.com)



114