






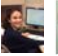





Research Activities in the Applied Genetics group

Peter M. Vallone, Ph.D.
 Leader, Applied Genetics Group NIST
 May 31, 2018
 University of Copenhagen








Applied Genetics Group – Forensic & Clinical Genetics

Topics

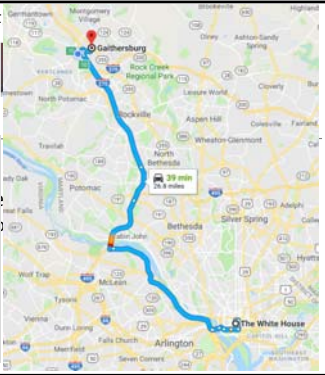
- NIST Standard Reference Materials (SRMs)
 - 2372a (Human DNA Quantitation), 2391d (PCR-based Profiling)
- STRBase 2.0
- Rapid DNA Interlaboratory Study
- **Ongoing Sequencing Projects**

Applied Genetics Group – F

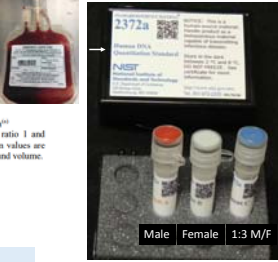
The National Institute of Standards and Technology (NIST) was founded in 1901 and is now part of the U.S. Department of Commerce. NIST is one of the nation's oldest physical science laboratories.

Congress established the agency to remove a major challenge to U.S. industrial competitiveness at the time—a second-rate measurement infrastructure that lagged behind the world.



SRM 2372a - Human DNA Quantitation Standard


- **On sale March 26, 2018**
- https://www-s.nist.gov/srmors/view_detail.cfm?srn=2372a
- Certified by digital PCR measurements



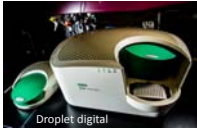
Component	Copy Number TM (per µL)	DNA TM (ng/µL)
A (red cap)	15.1 ± 1.5	49.8 ± 5.0
B (white cap)	17.5 ± 1.8	57.8 ± 5.8
C (blue cap)	14.5 ± 1.5	47.9 ± 4.8

To be used as a qPCR calibrant
 OR to assign a value to in house or commercial DNAs


Digital PCR Platforms at NIST




Fluidigm BioMark (2010)
Chamber digital



Bio-Rad QX100 (2012) upgraded to QX200 (2014)
Droplet digital



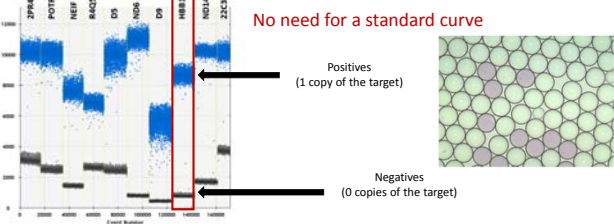
RainDance (2015)



Bio-Rad AutoDG (2015)

Digital PCR

Partitioning of DNA targets into individual chambers or droplets



dPCR is counting accessible amplifiable targets

Factors that affect quantification

$$\lambda = -\ln(N_{neg}/N_{tot})$$

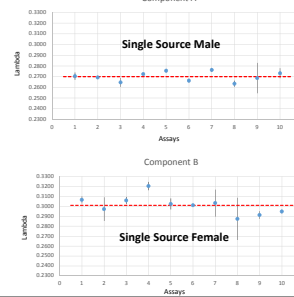
$$C = \lambda/(FV)$$

λ = number of targets per partition
 V = mean partition volume
 F = volume fraction of sample in the reaction mixture
 C = target concentration in a sample

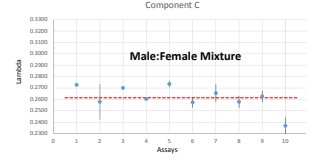
- Thermal cycler
- PCR assays – multiple targets
- Partition volume
- Treatment of artifacts
- Pipettes/balance
- Pre-treatment of DNA
- PCR master mix
- Primers/probes
- Lambda range

Not an exhaustive list

Quantification by dPCR a method to count DNA copies



Values assigned by 10 human custom genomic assays
Probing single copy targets in the human genome



All samples behave as expected with commercial qPCR assays (11 tested)

Converting copies per nanoliter to nanograms nuclear DNA per microliter

$$[nDNA]_{\mu L} = \left(\frac{\lambda \text{ copies of target}}{\text{droplet}} \right) \left(\frac{\mu L \text{ mixture}}{F \mu L \text{ sample}} \right) \left(\frac{\text{droplet}}{V \text{ mixture}} \right) \left(\frac{HHGE}{r \text{ target}} \right)$$

$$\left(\frac{n \text{ base pairs}}{HHGE} \right) \left(\frac{m \text{ g}}{\text{mol base pairs}} \right) \left(\frac{10^9 \text{ nL}}{6.022 \cdot 10^{23} \text{ base pairs}} \right) \left(\frac{10^9 \text{ nL}}{\mu L} \right) \left(\frac{10^9 \text{ ng}}{\text{g}} \right)$$

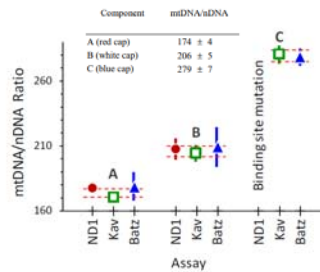
where r is the number of assay targets per human haploid genome equivalents (HHGE), n is the number of nucleotide base pairs (bp) per double-stranded HHGE, and m is the average molar mass of a bp in the DNA polymer.

For independent multiplicative factors such as these, the combined relative uncertainty of their product can be estimated from the square root of the sum-of-squares of the individual relative uncertainties [14, Section 5.1.6].

$$\frac{u([nDNA])}{[nDNA]} = \sqrt{\left(\frac{u(\lambda)}{\lambda} \right)^2 + \left(\frac{u(F)}{F} \right)^2 + \left(\frac{u(V)}{V} \right)^2 + \left(\frac{u(r)}{r} \right)^2 + \left(\frac{u(n)}{n} \right)^2 + \left(\frac{u(m)}{m} \right)^2}$$

Duewer DL, Kline MC, Ramsos EL, Toman B. Anal Bioanal Chem. 2018 410:2879-2887

SRM 2372a includes the ratio of mitochondrial to nuclear haploid genomes



mtDNA/nDNA ratio for three mitochondrial quantification assays optimized for dPCR.

SRM 2372a provides the ratio of mtDNA to gDNA, which bridges the gap between well characterized mtDNA quantification assays and availability of a commercial standard.

Candidate SRM 2391d PCR-based DNA profiling

- Successor to SRM 2391c
- Similar format – five tubes
 - A-C three single source components
 - D one mixture; approximately 3:1 (F:M)
 - E one component: cells spotted on FTA paper (from cell lines)
- Components A-D are DNA extracted from blood (not cell lines)
- Certified allele calls for core STR loci
- Characterized by CE- and NGS-based methods



Supports the FBI Quality Assurance Standards

Autosomal STR Markers

- Autosomal STR Markers
- ThermoFisher CE STR kits
- Promega CE STR kits
- Qiagen Investigator CE STR kits
- Illumina NGS kit
- ThermoFisher NGS kits
- Promega NGS kits
- CODIS 20/ESS 12

24 Certified Autosomal STR Markers
 1 Reference Autosomal STR Marker
 15 Information Autosomal STR Markers

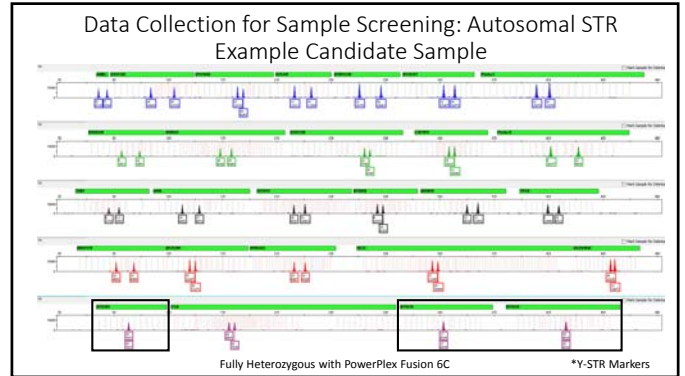
Autosomal STR Marker List	Applied Biosystems	Qiagen	ThermoFisher	Promega	Illumina	CODIS 20/ESS 12	Reference	Information
STR1								
STR2								
STR3								
STR4								
STR5								
STR6								
STR7								
STR8								
STR9								
STR10								
STR11								
STR12								
STR13								
STR14								
STR15								
STR16								
STR17								
STR18								
STR19								
STR20								
STR21								
STR22								
STR23								
STR24								
STR25								
STR26								
STR27								
STR28								
STR29								
STR30								
STR31								
STR32								
STR33								
STR34								
STR35								
STR36								
STR37								
STR38								
STR39								
STR40								
STR41								
STR42								
STR43								
STR44								
STR45								
STR46								
STR47								
STR48								
STR49								
STR50								

Y-STR Markers

Y-STR Markers
 ThermoFisher CE STR kits
 Promega CE STR kits
 Qiagen Investigator CE STR kits
 Illumina NGS kit
 ThermoFisher NGS kits
 Promega NGS kits

Y-STR Marker List	GoldenGate	Yfiler	Yfiler Plus	PPF Fusion 6C	PowerPlex Y25	Zenode Y25	Phosphor ID G2P	Phosphor ID Maxima ID G2P	PowerSeq 462Y	Genome Y1000	Genome Y1000	Information Y1000
DYS19									X			
DYS385a/b									X			
DYS389I									X			
DYS389II									X			
DYS391									X			
DYS392									X			
DYS437									X			
DYS438									X			
DYS439									X			
DYS448									X			
DYS459									X			
DYS460									X			
DYS466									X			
DYS468									X			
DYS481									X			
DYS505									X			
DYS518									X			
DYS522									X			
DYS523									X			
DYS524									X			
DYS525									X			
DYS526									X			
DYS527									X			
DYS528									X			
DYS529									X			
DYS530									X			
DYS531									X			
DYS532									X			
DYS533									X			
DYS534									X			
DYS535									X			
DYS543									X			
Y-DATA-H4									X			
DYS392b1									X			

23 Certified Y-STR Markers
0 Reference Y-STR Markers
6 Information Y-STR Markers



Data Collection for Sample Screening: Y-STR

Yfiler Plus Profile

YHRD: No matches in 188,209 Haplotypes (using Minimal Haplotype) <https://yhrd.org>

Whit Athey's Haplogroup Predictor: E1b1a <http://www.hprg.com/hapest5/hapest5a/hapest5.htm?order=num>

Haplogroup	Fitness score	Probability (%)
E1b1a	10	100
E1b1b	18	0.0
G2a	20	0.0
G2c	3	0.0
H	25	0.0
I1	3	0.0
I2a	20	0.0
I2a1	3	0.0
I2b	7	0.0
I2b1	14	0.0
I1	11	0.0
I2a1b	3	0.0
I2a1b	4	0.0
I2a1 x I2a1-Mb	11	0.0
I2b	9	0.0
L	13	0.0
N	2	0.0
O	17	0.0
R1a	11	0.0
R1b	4	0.0
T	16	0.0

Information for additional marker systems

Support the adoption of new markers and technology platforms

- Mitochondrial genome sequence
- Identity SNPs – for degraded samples
- Ancestry SNPs – biogeographical ancestry prediction
- Phenotype SNPs – eye and hair color prediction

Data Collection for Sample Screening: SNPs

ForenSeq SNP Phenotype and Ancestry Estimation

Brown	0.16
Red	0.00
Black	0.84
Blond	0.00

Intermediate	0.00
Green	1.00
Blue	0.00

Distance to Nearest Centroid	3.36
------------------------------	------

Additional markers to be characterized:
 X-STRs, Indels, INNULS, other SNP Panels, and Microhaplotypes

Data Collection for Sample Screening: mtDNA

Illumina mtDNA Whole Genome Sequencing protocol with Nextera XT Sample Prep Kit

L1b1a12: 15 haplotypes

EMPOP results: https://empop.online/haplotypes#matches_details

Haplogroup	Ancestry	Match
L1b1a12	African	unique

STRBase 2.0 Draft image

STRBase 2.0

- First round of development
 - STR fact sheets (for 24 loci)
 - Variant allele reporting
 - Goal is to have a beta site up this summer
- Provide search, sort, and download functionalities
- Automated submission of variant alleles
- Embedded viewer for STR sequence and presentations
- Other ideas or functionalities – let me know

STR Fact sheet: example D1S1656

Visualization of the Sequence

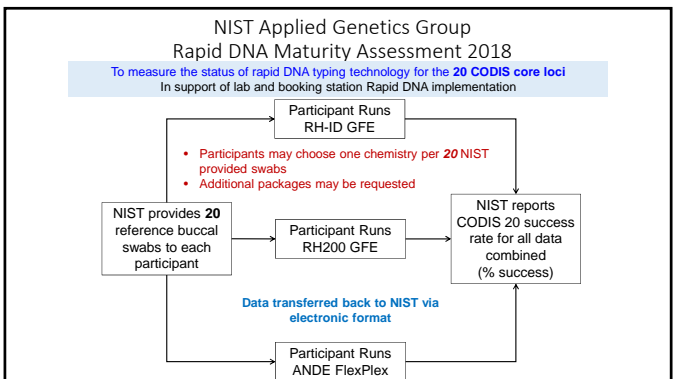
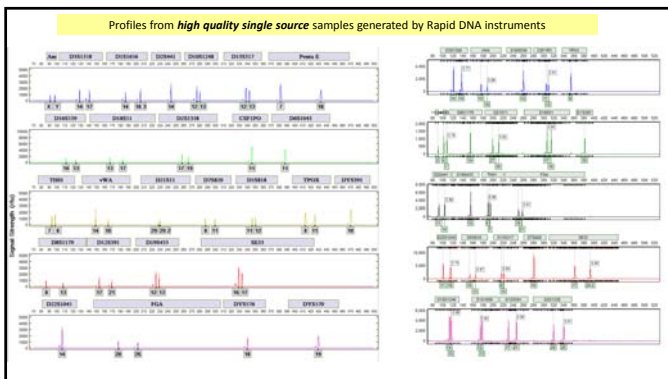
Identify surrounding sequence and provide observed SNPs

Viewer will link to NCBI Bioproject

Rapid DNA Assessment III


- Summer 2018
- Core **20** STR markers
- **Projected: 8 labs and 2 vendors**
- ANDE and IXI platforms (**new kits, configurations**)
- **20** samples per lab (single source swab)
- Currently collecting single source swabs

Supporting the use of Rapid DNA in the booking station



Sequencing Projects

- FGx and S5 platforms
- **1036 population samples**
- **Highly polymorphic locus SE33**
- **STRSeq resource**
- **Nomenclature support**
- **Sensitivity studies**
- Concordance projects



Applications for Sequencing STRs

- Targeted sequencing of STRs
 - STR motif sequence variation; flanking region variation
 - Further understand simple versus complex repeat motifs
 - Characterize stutter
- Applications
 - One to one matching?
 - With the new U.S. core loci we are already quite high (>10⁻²⁰)
 - Partial profiles
 - Kinship
- Greater degree of multiplexing
 - Not confined by dye colors; smaller PCR amplicons (degraded samples)
 - PCR for sample enrichment
 - **Still using PCR** – stochastic effects, stutter
- Mixtures
 - Resolve alleles identical by length, but differ by sequence
 - Separate stutter from low level contributors (based on sequence)
 - A sequenced allele *may* have a lower frequency (higher RMP or LR)

Allele frequencies of sequenced STR alleles are needed to formally apply this gain in information -> **Generate population data!**

Sequencing Forensic STRs in Population Samples

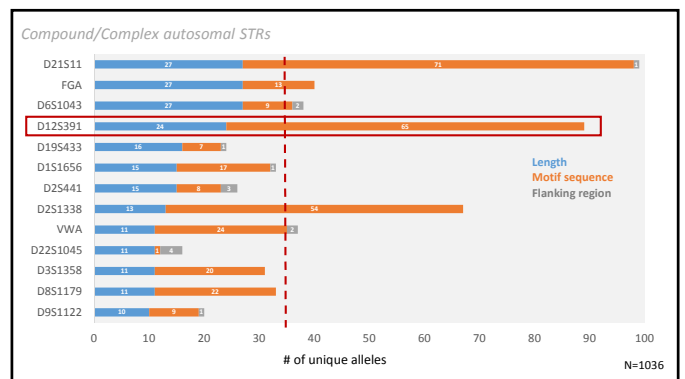
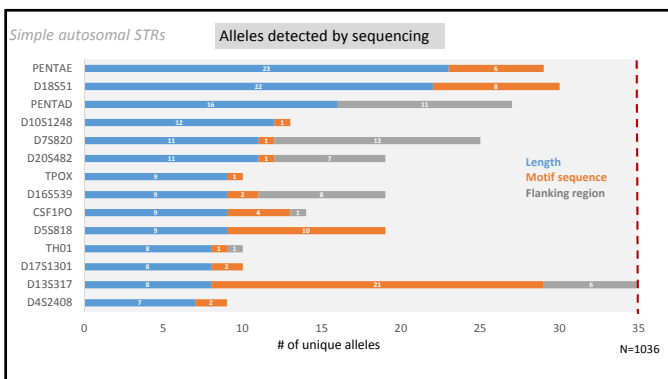
When a match is made in a forensic case, allele frequencies are used to calculate how common or rare the DNA profile is in a given population

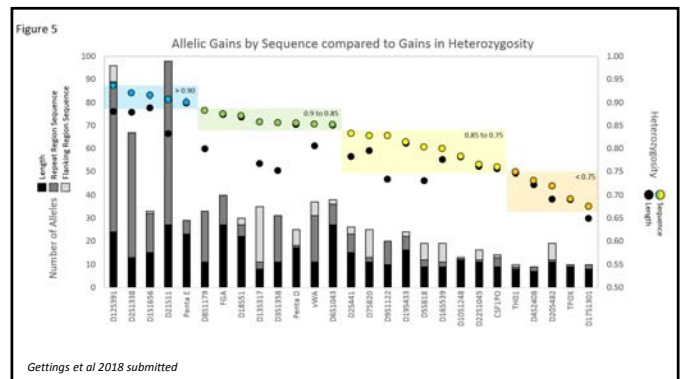
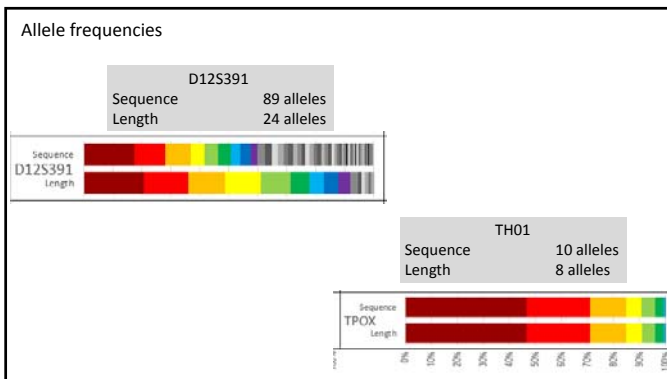
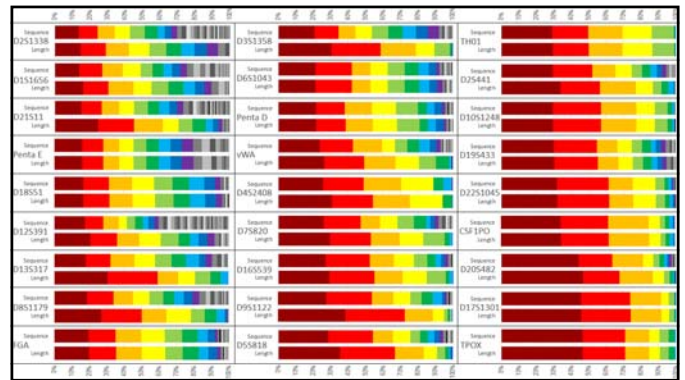
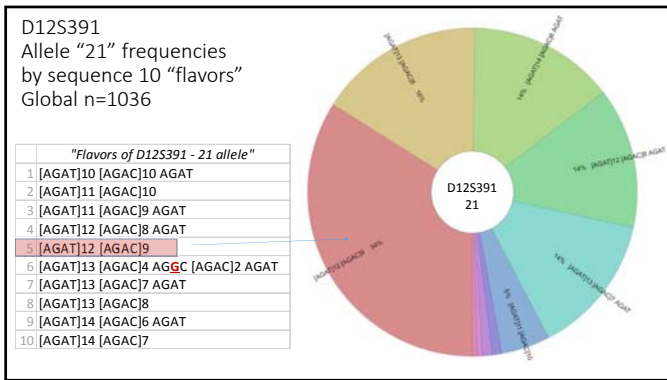
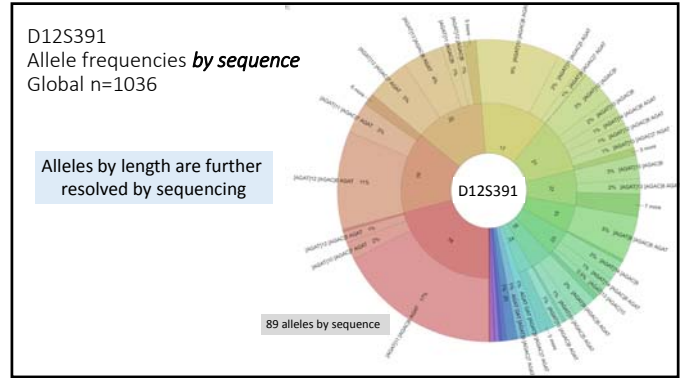
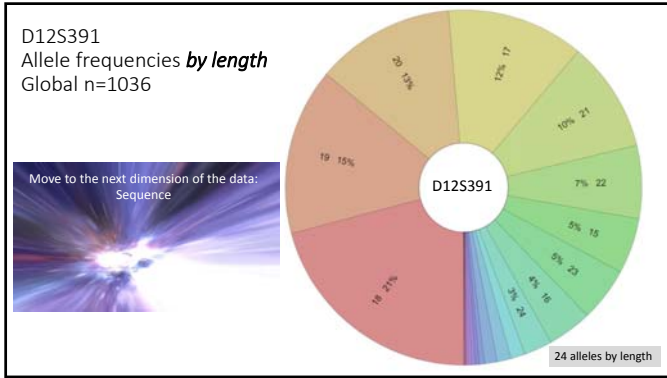
Example of **length** versus **sequence**-based frequency calculation:

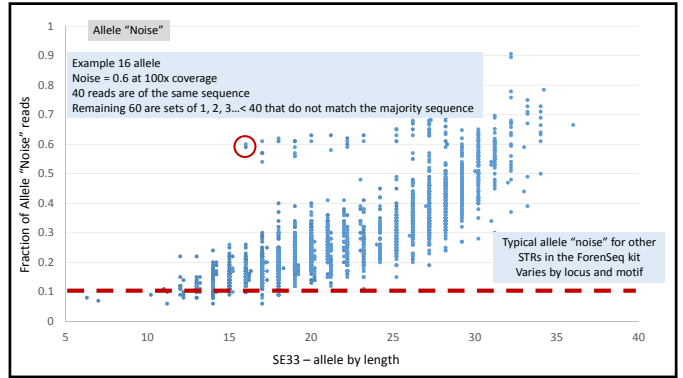
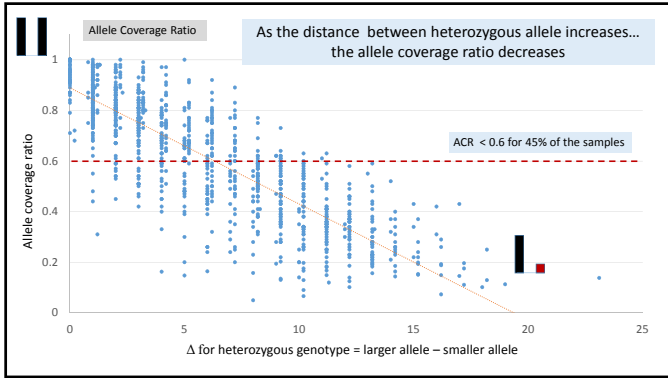
D4S2408						Length	Sequence
Allele	N	Freq	Sequence	Allele	N	Freq	
7	1	0.6%	[ATCT]7		1	0.6%	8,9 [ATCT]8, [ATCT]9
8	23	14.4%	[ATCT]8		23	14.4%	2pq 2pq
9	60	37.5%	[ATCT]9		18	11.3%	2*0.144*0.375 2*0.144*0.113
			[ATCT] G TCT [ATCT]7	42	26.3%		
10	53	33.1%	[ATCT]10		53	33.1%	0.108 0.033
11	21	13.1%	[ATCT]11		21	13.1%	
12	2	1.3%	[ATCT]12		2	1.3%	1 in 9.3 1 in 30.7

Sequencing of 1036 NIST population samples

- Work performed on Illumina FGx – ForenSeq kit
- Allele calls were made with Illumina-UAS and STRait Razor and compared to CE length-based calls (**high confidence**)
- Will include flanking region variation (SNPs, indels)
- Purpose: provide sequence allele frequencies for four U.S. Population groups
 - U.S.: Caucasian, African American, Hispanic, Asian
- The manuscript was submitted in **late May**
 - Focus on the autosomal loci







STRSeq

Forensic Science International: Genetics

Research paper
STRSeq: A catalog of sequence diversity at human identification Short Tandem Repeat loci

Katherine Butler Gettings¹, Lisa A. Borsak², David Ballard³, Martin Budser⁴, Bruce Budwala⁵, Laurence Devener⁶, Jonathan King⁷, Wilbur Payne⁸, Christopher Phillips⁹, Peter M. Vallone¹⁰

- Bioproject hosted at NCBI to catalog **unique** STR alleles
 - <https://www.ncbi.nlm.nih.gov/bioproject/380127>
- Fully annotated sequence of the STR amplicon
- Each STR allele will have an accession number

Resource Name	Number of Links
Sequence data	
Nucleotide (Genomic DNA)	1145
Publications	2

<https://strider.online/>

Collaborative Effort (NIST, KCL, UNT, USC)

1786 + 1043 + 839 + 944 = 4612 samples

Example given for D12S391

(a) Venn diagram showing the overlap of unique alleles between NIST (105), UNT (80), USC (96), and KCL (97).

(b) Substitution strategy for 157 unique sequence based alleles observed in the D12S391 locus. The 157 unique alleles generated at NIST from the basis of 5766 records, subsequent observations from KCL, UNT, and USC will add records for segments generated in each laboratory for which records do not already exist (98, 9, and 18 records, respectively).

- Initial NIST Data = 105 Records
- +25 KCL Alleles = 130 Records
- +9 UNT Alleles = 139 Records
- +18 USC Alleles = 157 Records

STRSeq

- A unique record for each observed allele
- CE-based allele call
- Genomic coordinates
- Platform/kit used to generate sequence

- Bioproject hosted at NCBI to catalog
 - <https://www.ncbi.nlm.nih.gov/bioproject/380127>
- Fully annotated sequence of the:
- Each STR allele will have an acces

Items: 16

1. 82 bp linear DNA
Accession: M098827.1 | GI: 139120738
BioProject: EMM04 | Sequence: GenBank | FASTA | GenScript
2. 167 bp linear DNA
Accession: M098827.1 | GI: 139120737
BioProject: EMM04 | Sequence: GenBank | FASTA | GenScript
3. 51 bp linear DNA
Accession: M098827.1 | GI: 139120736
BioProject: EMM04 | Sequence: GenBank | FASTA | GenScript
4. 167 bp linear DNA
Accession: M098827.1 | GI: 139120735
BioProject: EMM04 | Sequence: GenBank | FASTA | GenScript
5. 155 bp linear DNA
Accession: M098827.1 | GI: 139120734
BioProject: EMM04 | Sequence: GenBank | FASTA | GenScript

Nomenclature Support

Thoughts on reporting?

STR allele sequence variation: Current knowledge and future issues

Forensic Science International: Genetics

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

<https://strider.online>

Defining annotated STR reference sequences

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

C. Phillips¹, K. Butler Gettings², J.L. King³, D. Ballard⁴, M. Budser⁵, L. Borsak⁶, W. Payne⁷

NGS – MPS – Sequencing noise

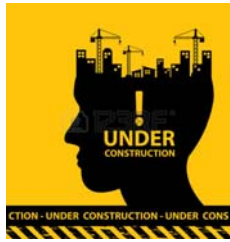
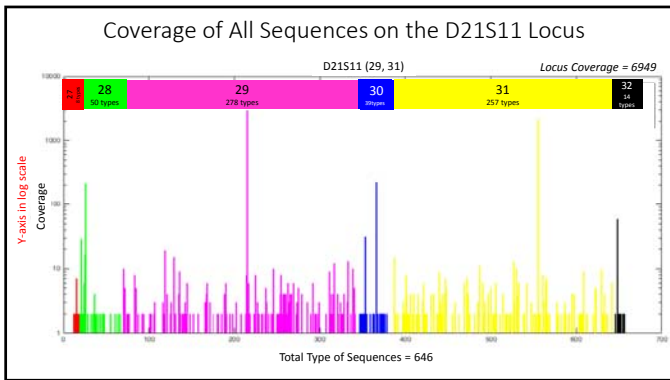
Thoughts

Sarah Riman
Hari Iyer (NIST statistical engineering division)


- How should sequence 'noise' be characterized?
- Coverage or normalized coverage?
- Evaluate locus specific performance
- Setting thresholds?
 - Need to understand data first
- True allele versus noise
 - for now: including artifacts and stutter as 'noise'
- Zygosity

First steps

- Understand the data from sensitivity studies
- Develop tools to assess

Noise Thresholds for NGS Data



Understand the characteristics of single source DNA profiles generated by the NGS system by evaluating...

- Receiver Operating Characteristic (ROC) curves to define where alleles can be clearly separated from noise (attributed to either stutter or random causes)
- Zygosity to minimize the risks of misidentifying a heterozygote as a homozygote or a homozygote as heterozygote

In this study we evaluated reference samples ONLY
We did not test any unknown (crime-stain), mixtures or degraded samples.

Experiments

- Promega PowerSeq Auto 46GY
- Illumina MiSeq (v3 and 92 samples/run)
- Library preparation; TruSeq and Kapa
- Normalized and non-normalized

Sample	500	250	125	60	30	15
Sample 1	500	250	125	60	30	15
Sample 2	500	250	125	60	30	15
Sample 3	500	250	125	60	30	15

Single source samples – amplified in triplicate

Locus Name	Size (bp)	Locus Name	Size (bp)
D1S1179	203-206	D1S281	167-170
D1S1179	203-206	D1S291	188-194
D1S1179	203-206	D1S301	202-203
D1S1179	203-206	D1S311	204-204
D1S1179	203-206	D1S321	204-204
D1S1179	203-206	D1S331	204-204
D1S1179	203-206	D1S341	204-204
D1S1179	203-206	D1S351	204-204
D1S1179	203-206	D1S361	204-204
D1S1179	203-206	D1S371	204-204
D1S1179	203-206	D1S381	204-204
D1S1179	203-206	D1S391	204-204
D1S1179	203-206	D1S401	204-204
D1S1179	203-206	D1S411	204-204
D1S1179	203-206	D1S421	204-204
D1S1179	203-206	D1S431	204-204
D1S1179	203-206	D1S441	204-204
D1S1179	203-206	D1S451	204-204
D1S1179	203-206	D1S461	204-204
D1S1179	203-206	D1S471	204-204
D1S1179	203-206	D1S481	204-204
D1S1179	203-206	D1S491	204-204
D1S1179	203-206	D1S501	204-204
D1S1179	203-206	D1S511	204-204
D1S1179	203-206	D1S521	204-204
D1S1179	203-206	D1S531	204-204
D1S1179	203-206	D1S541	204-204
D1S1179	203-206	D1S551	204-204
D1S1179	203-206	D1S561	204-204
D1S1179	203-206	D1S571	204-204
D1S1179	203-206	D1S581	204-204
D1S1179	203-206	D1S591	204-204
D1S1179	203-206	D1S601	204-204
D1S1179	203-206	D1S611	204-204
D1S1179	203-206	D1S621	204-204
D1S1179	203-206	D1S631	204-204
D1S1179	203-206	D1S641	204-204
D1S1179	203-206	D1S651	204-204
D1S1179	203-206	D1S661	204-204
D1S1179	203-206	D1S671	204-204
D1S1179	203-206	D1S681	204-204
D1S1179	203-206	D1S691	204-204
D1S1179	203-206	D1S701	204-204
D1S1179	203-206	D1S711	204-204
D1S1179	203-206	D1S721	204-204
D1S1179	203-206	D1S731	204-204
D1S1179	203-206	D1S741	204-204
D1S1179	203-206	D1S751	204-204
D1S1179	203-206	D1S761	204-204
D1S1179	203-206	D1S771	204-204
D1S1179	203-206	D1S781	204-204
D1S1179	203-206	D1S791	204-204
D1S1179	203-206	D1S801	204-204
D1S1179	203-206	D1S811	204-204
D1S1179	203-206	D1S821	204-204
D1S1179	203-206	D1S831	204-204
D1S1179	203-206	D1S841	204-204
D1S1179	203-206	D1S851	204-204
D1S1179	203-206	D1S861	204-204
D1S1179	203-206	D1S871	204-204
D1S1179	203-206	D1S881	204-204
D1S1179	203-206	D1S891	204-204
D1S1179	203-206	D1S901	204-204
D1S1179	203-206	D1S911	204-204
D1S1179	203-206	D1S921	204-204

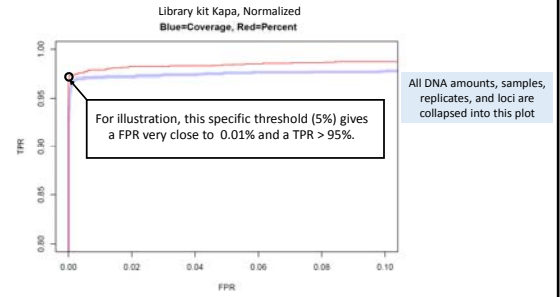
ROC-defined analytical threshold

- Is a two dimensional chart which plot the **true positive** versus the **false positive** rates for a given parameter
- Is performed to determine which AT would lead to optimal levels of detection where error rates are minimized
- The true positive rate represents the proportion of true allele sequences known to be present at a specific locus/method/DNA amount
- The false positive rate generated at a locus of interest/method/DNA amount represents the proportion of noise sequences falsely classified as true sequences
- **Assumption of a certain distribution is not required with ROCs**
- Measures performance of different ATs and gives proportion of false positives and negatives.

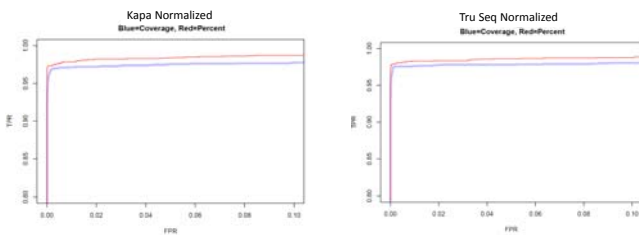
ROCs can be used to determine:

- ❖ A global analytical threshold
- ❖ A DNA amount-dependent threshold
- ❖ A DNA amount and locus-dependent threshold

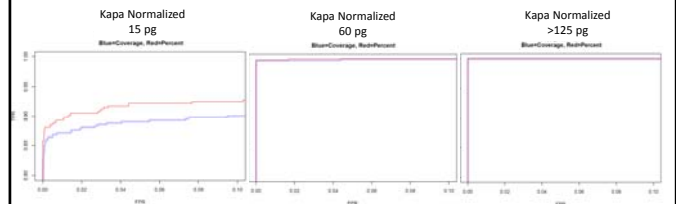
Comparison of Coverage versus Percent ROC



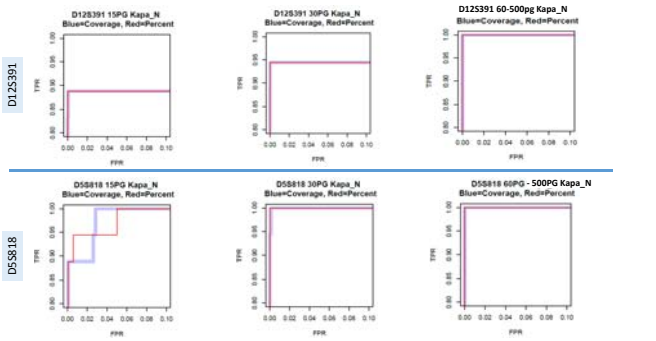
ROCs as a function of library preparation kit



ROCs as a function of DNA input amount



A DNA amount and locus-dependent threshold (all Kapa Normalized)

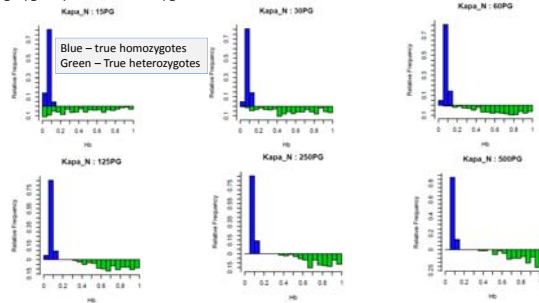


Inferring Zygosity from Heterozygote Balance

Two primary risks associated with the process of allele designation:

- A heterozygote is inaccurately assigned as a homozygote due to:
 - Allele drop-out has occurred
 - Heterozygote imbalance has resulted in one of the alleles being interpreted as a stutter
- A homozygote is inaccurately called as a heterozygote due to:
 - A large stutter band is within range to be designated as an allele
 - A drop-in event occurs

Inferring Zygosity from Heterozygote Balance



A comparison of the distribution of the homozygotes and heterozygotes showed marked differences associated with the differences in the DNA input.

Further work

- Incorporate stutter and accountable artifacts into the ROCs
- Perform further sensitivity experiments as needed
- Create mock casework type samples/mixtures
 - Derive and test thresholds

Thank you for your attention! Questions?

Contact: Peter.Vallone@nist.gov



- **Funding**
 - NIST Special Programs Office: *Forensic DNA*
 - FBI Biometrics Center of Excellence: *Forensic DNA Typing as a Biometric tool.*
 - NIJ: *STRSeq and Nomenclature*
 - DHS S&T: *Rapid DNA for Kinship*
- **Disclaimer** - Points of view in this document are those of the authors and do not necessarily represent the official position or policies of the U.S. Department of Commerce or the 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 NIST, nor does it imply that any of the materials, instruments, or equipment identified are necessarily the best available for the purpose.
- All work presented has been reviewed and approved by the NIST Human Subjects Protections Office.