



# Amplicon preparation and data analysis for autosomal STR profiling using next generation sequencing

In Seok Yang · Eun Hye Kim · Kyoung-Jin Shin  
Department of Forensic Medicine, Yonsei University College of Medicine, Seoul, Korea

## Introduction

STR markers are highly variable among individuals, thereby making these STRs effective for human identification and kinship testing. Length-based genotyping using capillary electrophoresis (CE) has been the standard approach to forensic STR analysis. Recently emerged next generation sequencing (NGS) has the potential to be ultimate genotyping platform for forensic purpose. Still, amplicon sequencing method for STR profiling using NGS and its resultant data analysis protocol were not well established. Therefore, we tested 3 different amplicon preparation methods for NGS library and sought to present appropriate NGS data analysis protocol for PowerPlex 16 STRs profiling in this study.

## Materials and Methods

### 1. STR loci and primer design

Dataset 1 : PCR product using PowerPlex 16 HS system

Dataset 2: PCR product using NGS adaptor fusion primers

Forward primer (Primer A-Key)

5'-CGTATCGCCTCCCTCGCCGATCAG-(template-specific-sequence)-3'

Reverse primer (Primer B-Key)

5'-CTATGCGCCTTGCACGCCGCTCAG-(template-specific-sequence)-3'

Dataset 3: PCR product using original primers of PowerPlex 16 system

	Type 1 primer set Dataset 1	Type 2 primer set Dataset 2	Type 3 primer set Dataset 3
Fluorescence dye attached	+	-	-
Adaptor sequence included	-	+	-

### 2. STR amplicon preparation

STR amplicons were prepared using PowerPlex 16 protocol with 34 thermal cycles from 1 ng of 2800M Control DNA sample. Then the amplicons were purified using Qiaquick PCR Purification Kit.

### 3. NGS library preparation

For PCR products of dataset 1 and 3, NGS libraries were prepared by ligation of NGS adaptors.

For PCR products of dataset 2, NGS library preparation was omitted, since PCR amplification was carried out with NGS adaptor fusion primer.

Size selection using AMPure beads was subsequently performed to remove small fragments.

### 4. NGS data generation

NGS datasets were obtained by emulsion PCR of the NGS libraries and subsequent sequencing on 1/8 of a GS Pico TiterPlate using Roche (454) GS FLX Titanium Sequencer.

### 5. NGS data analysis

NGS data analysis in this study basically follows the protocol presented by Borman et al. (Biotechniques, 2012) with modifications on coverage calculation and STR allele assignment.

STR reference sequences were built using known STR repeat structure in STRbase (<http://www.csl.nist.gov/strbase/>). 5' and 3' flanking sequences were from human reference sequence GRCh37/hg19.

Sequence reads of 3 NGS datasets were aligned onto the STR reference sequence using Bowtie 2 program on Linux operating system. SAMtools and BEDTools were subsequently used for format conversion of alignment outputs.

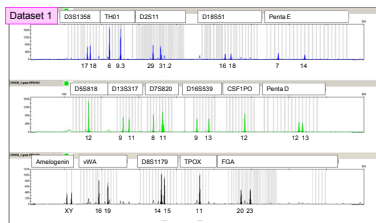
For coverage calculation, aligned sequence reads were filtered based on breadth and mapping quality of them.

### 6. Assignment of STR alleles

STR alleles were assigned by using coverage threshold of 20%, which were empirically determined for single-source samples.

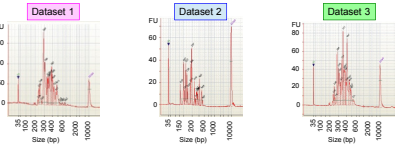
## Results

### 1. CE profile of STR amplicon



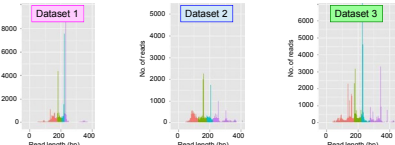
Full STR genotypes for 16 loci could be obtained from all STR amplicons that were generated with 3 different primer sets.

### 2. NGS library profiles



All NGS libraries were ranged between 150 and 600 bp. These library profiles were formed in relatively smaller size range than those of nebulized DNA fragments for sequencing of whole genome or mitochondrial DNA.

### 3. Read length distributions



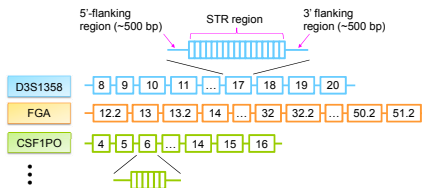
Through duplicate experiment, three different NGS datasets could be obtained consistently. Different patterns on read length distribution were observed among 3 NGS datasets.

### 4. Total number of sequence reads

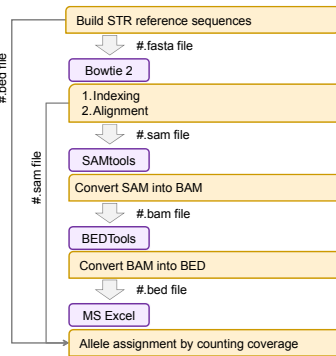
	Dataset 1	Dataset 2	Dataset 3
Data quantity (total number of sequence reads)	143,404	82,371	155,781
	***	*	***

NGS datasets obtained by PCR amplification f target STR loci followed by adaptor ligation (datasets 1 and 3) showed higher data quantity than those from PCR products generated with NGS adaptor fusion primers (dataset 2).

### 5. Schematic view of STR reference sequence



### 6. STR profiling protocol from NGS data



### 7. Alignment results of NGS datasets

STR	Dataset 1		Dataset 2		Dataset 3	
	Number of aligned reads	Reads (coverage)	Number of aligned reads	Reads (coverage)	Number of aligned reads	Reads (coverage)
D3S1358	1135	1210	1145	1680	1272	1520
D5S818	1548	875	1182	1274	808	4446
D7S820	8540	1869	1669	2206	1871	5230
D13S317	60	36	4844	2953	1991	11346
D16S11	16739	3221	10509	6337	1722	10020
D18S21	87	46	2775	1157	1138	11392
D19S11	509	113	146	2738	1661	1077
D21S11	3999	1203	20206	12134	1245	3885
CSF1PO	112	81	12	7669	1822	1757
FGA	305	81	122	6505	1278	1327
TH01	277	84	186	867	425	482
TH02	82	81	151	868	540	705
TH03	17260	1727	10133	13023	1541	7870
TH04	12	15	15	2345	1266	877
TH05	85	35	25	5305	1463	1817
Total	140394	44077	76117	76680	39142	27526

The number of sequence reads mapping to each STR locus appeared differently among 3 NGS datasets.

### 8. Coverage of STR alleles

STR	Start	End	Allele	Dataset 1			Dataset 2			Dataset 3		
				Count	Coverage	Allele	Count	Coverage	Allele	Count	Coverage	Allele
CSF1PO	1536	1537	CSF1PO_5	0	0	0	0	0	0	0	0	
	1536	1539	CSF1PO_7	0	0	0	0	0	0	0	0	
	1536	1541	CSF1PO_9	0	0	0	0	0	0	0	0	
	1536	1543	CSF1PO_11	0	0	0	0	0	0	0	0	
	1536	1545	CSF1PO_13	0	0	0	0	0	0	0	0	
	1536	1547	CSF1PO_15	0	0	0	0	0	0	0	0	
	1536	1549	CSF1PO_17	0	0	0	0	0	0	0	0	
	1536	1551	CSF1PO_19	0	0	0	0	0	0	0	0	
	1536	1553	CSF1PO_21	42	100	1180	42	100	1180	42	100	1180
	1536	1555	CSF1PO_23	78	180	2180	78	180	2180	78	180	2180
D18S11	1014	1025	D18S11_1	0	0	0	0	0	0	0	0	
	1014	1027	D18S11_3	0	0	0	0	0	0	0	0	
	1014	1029	D18S11_5	0	0	0	0	0	0	0	0	
	1014	1031	D18S11_7	0	0	0	0	0	0	0	0	
	1014	1033	D18S11_9	0	0	0	0	0	0	0	0	
	1014	1035	D18S11_11	0	0	0	0	0	0	0	0	
	1014	1037	D18S11_13	0	0	0	0	0	0	0	0	
	1014	1039	D18S11_15	0	0	0	0	0	0	0	0	
	1014	1041	D18S11_17	0	0	0	0	0	0	0	0	
	1014	1043	D18S11_19	0	0	0	0	0	0	0	0	

NGS datasets in high quality order were dataset 3, dataset 2, and dataset 1. Thus, it revealed that NGS dataset with the highest quality was obtained by PCR amplification using original primers of PowerPlex 16 system (no fluorescence dye and NGS adaptor).

### 9. An example of STR allele assignment

STR	Start	End	Allele	Coverage	Total coverage	Assignment
D18S11	1014	1025	D18S11_1	0	1844	0/1844
D18S11	1014	1027	D18S11_3	272	1844	14.75%
D18S11	1014	1029	D18S11_5	858	1844	46.53%
D18S11	1014	1031	D18S11_7	0	1844	0%
D18S11	1014	1033	D18S11_9	873	1844	47.34%
D18S11	1014	1035	D18S11_11	119	1844	6.48%
D18S11	1014	1037	D18S11_13	119	1844	6.48%
D18S11	1014	1039	D18S11_15	0	1844	0%
D18S11	1014	1041	D18S11_17	0	1844	0%
D18S11	1014	1043	D18S11_19	0	1844	0%

Most STR alleles could be assigned exactly with coverage threshold of 20%.

### 10. Accuracy of STR allele assignment

STR	# results	Dataset 1		Dataset 2		Dataset 3	
		Allele	Rate	Allele	Rate	Allele	Rate
D3S1358	12	12	100%	12	100%	12	100%
D5S818	12	12	100%	12	100%	12	100%
D7S820	11	11	100%	11	100%	11	100%
D13S317	11	11	100%	11	100%	11	100%
D16S11	11	11	100%	11	100%	11	100%
D18S21	11	11	100%	11	100%	11	100%
D19S11	11	11	100%	11	100%	11	100%
D21S11	11	11	100%	11	100%	11	100%
CSF1PO	11	11	100%	11	100%	11	100%
FGA	11	11	100%	11	100%	11	100%
TH01	11	11	100%	11	100%	11	100%
TH02	11	11	100%	11	100%	11	100%
TH03	11	11	100%	11	100%	11	100%
TH04	11	11	100%	11	100%	11	100%
TH05	11	11	100%	11	100%	11	100%
Total	140394	140394	100%	140394	100%	140394	100%

Accuracy of STR allele assignment: Dataset 1 (★★), Dataset 2 (★), Dataset 3 (★★★). NGS datasets in high accuracy order on STR allele assignment were dataset 3, dataset 1, and dataset 2. Thus, it revealed that NGS dataset with accurate STR allele call was obtained by PCR amplification using original primers of PowerPlex 16 system (no fluorescence dye and NGS adaptor).

## Conclusion

Through comparison among three NGS datasets, it revealed that NGS dataset with the highest quality for STR profiling was obtained by amplicon preparation with template specific primer followed by adaptor ligation. Therefore, the amplicon preparation method and NGS data analysis protocol presented from this study will be helpful to interpret and analyze STR profile for forensic purpose using a NGS platform.

## Acknowledgments

This work was supported by the research project for practical use and advancement of forensic DNA analysis of Supreme Prosecutors' Office, Republic of Korea (1333-304, 2012).