

## Introduction

Because of its capability to produce massively parallel sequencing data and to provide sequence variations, the prospect of NGS in STR typing has been explored by many forensic laboratories using in-house multiplex PCR systems or commercial kits. Meanwhile, the expanded CODIS Core Loci were newly announced to increase the international data compatibility and to maximize the power of forensic DNA databases. Accordingly, we developed a multiplex PCR system for the NGS analysis of 25 forensic markers which include 20 expanded CODIS Core Loci (D1S1656, D2S441, D2S1338, D3S1358, D5S818, CSF1PO, D7S820, D8S1179, D10S1248, D12S391, D13S317, TPOX, D16S539, D18S51, D19S433, D21S11, D22S1045, FGA, TH01 and vWA) and additional 5 Loci (Penta D, Penta E, D6S1043, DYS391 and amelogenin). To investigate the sequence variations of 24 global STR loci, NGS data of 160 unrelated Koreans were analyzed. In addition, we present the observed sequence variations with allele frequencies and forensic parameters in Koreans.

## Materials and Methods

### 1. DNA samples

DNA was extracted from buccal swab samples of 160 unrelated Koreans using a QIAamp DNA Mini Kit (Qiagen) and quantified using a NanoDrop<sup>®</sup> ND-1000 Spectrophotometer (NanoDrop Technologies). Finally diluted 1 ng/ul of DNA was prepared and used. The study was approved by the Institutional Review Board of Severance Hospital, Yonsei University in Seoul, Korea.

### 2. Construction of multiplex PCR system

The STR markers selected for a multiplex PCR were composed of 25 forensic markers that included 20 expanded CODIS Core Loci and additional 5 Loci (Penta D, Penta E, D6S1043, DYS391 and amelogenin) shown in Fig. 1. Primers for PCR amplification were designed using the Primer3 (<http://bioinfo.ut.ee/primer3-0.4.0/primer3/>) program such that the amplicon sizes of the 25 targeted markers were less than 250 bp, and the primers did not bind to the region with a greater than 1% mutation rate based on the NCBI SNP information (<http://www.ncbi.nlm.nih.gov/SNP/>).

### 3. Preparation of NGS library

We conducted two-step PCR amplifications to generate a library using primers with a modification referring to the sequence information of Nextera<sup>®</sup> system (Illumina). The first PCR targeted the autosomal STR itself, and primer sequences included STR-specific sequences and read sequences. A second PCR was performed to add indices and platform-specific sequences. See detailed information on poster No. 79.

1) The first PCR - multiplex PCR was performed with 29 thermal cycles from each 1 ng of the sample and appropriate concentration of primers.

2) The second PCR - indexing PCR was performed with 15 thermal cycles from each 1.0 µl of 100-fold diluted the first PCR amplicons and Nextera XT Index Kit (Illumina).

Following PCR cleanup with 1.2× Agencourt<sup>®</sup> AMPure<sup>®</sup> XP beads (Beckman Coulter), the libraries were quantified using KAPA library quantification kits (KAPA Biosystems) and Agilent 2100 Bioanalyzer.

### 4. NGS run and data analysis

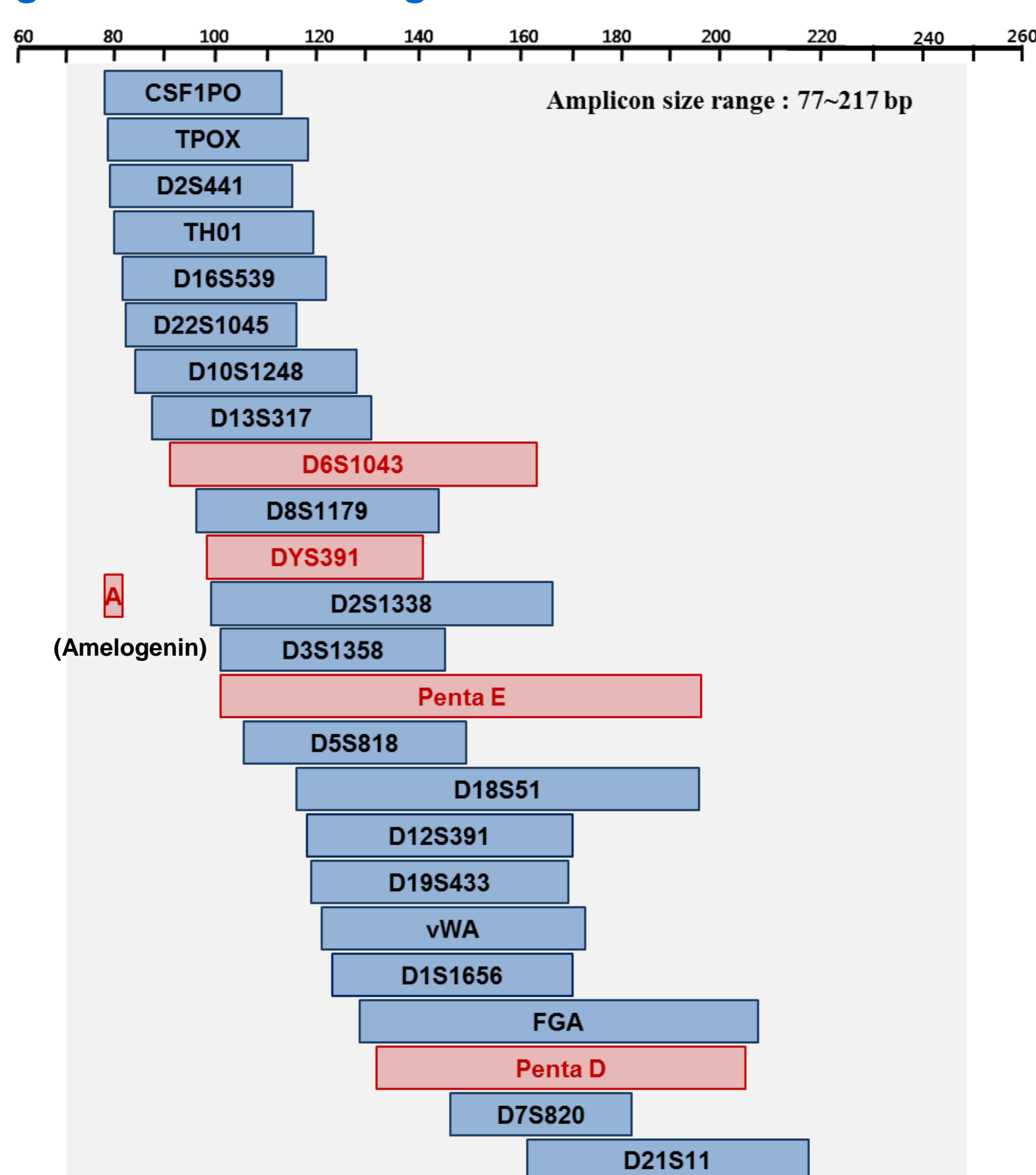
The barcoded libraries were normalized to 10nM and then pooled in equal volumes. Finally, the pooled library was sequenced on MiSeq<sup>™</sup> (Illumina) using a MiSeq Reagent Kit v2, 2x250 bp (Illumina). NGS data analysis basically follows the protocol presented by Bornman et al (Biotechniques, 2012). The process of NGS data analysis used in this study was illustrated in Fig. 2. The STR profiles obtained by two CE methods - Kplex-23 and Euplex-13 system (BioQuest) - were used as reference data for comparing the STR typing results from NGS.

### 5. Statistical analysis

Statistical parameters, including the value of the allele frequencies and observed heterozygosity ( $H_{obs}$ ), were calculated using counting methods. Forensic efficiency information was assessed by calculating expected heterozygosity ( $H_{exp}$ ), match probability (MP), power of discrimination (PD), polymorphic information content (PIC) and power of exclusion (PE).

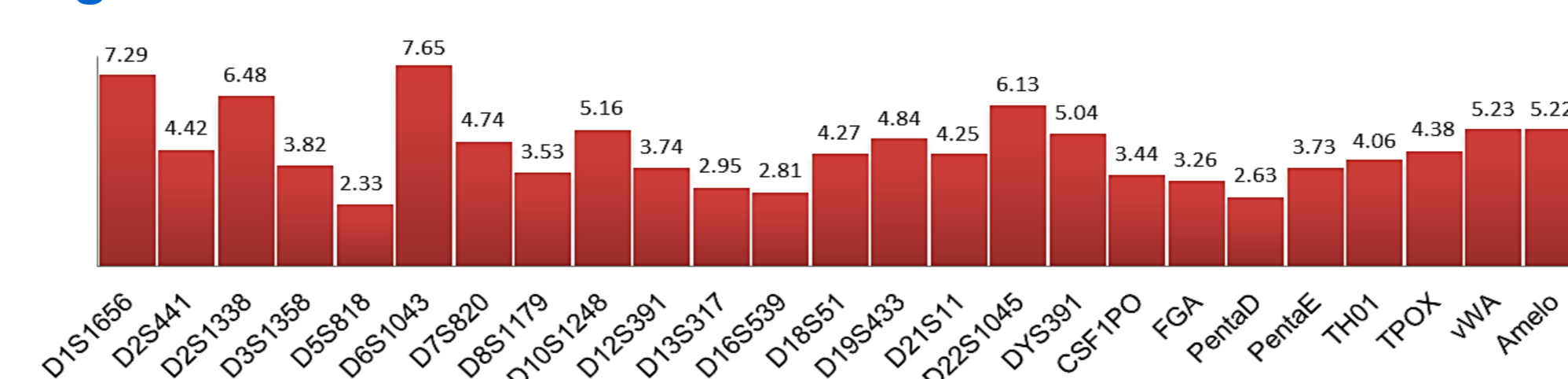
## Results

Fig. 1. Allelic size range for 25 forensic markers



Autosomal 20 expanded CODIS Core Loci were marked in blue boxes and additional 5 loci were marked in red boxes. Most markers were less than 200 bp. Those of Penta D and FGA were less than 210 bp. D21S11 showed a maximum size range less than 217 bp.

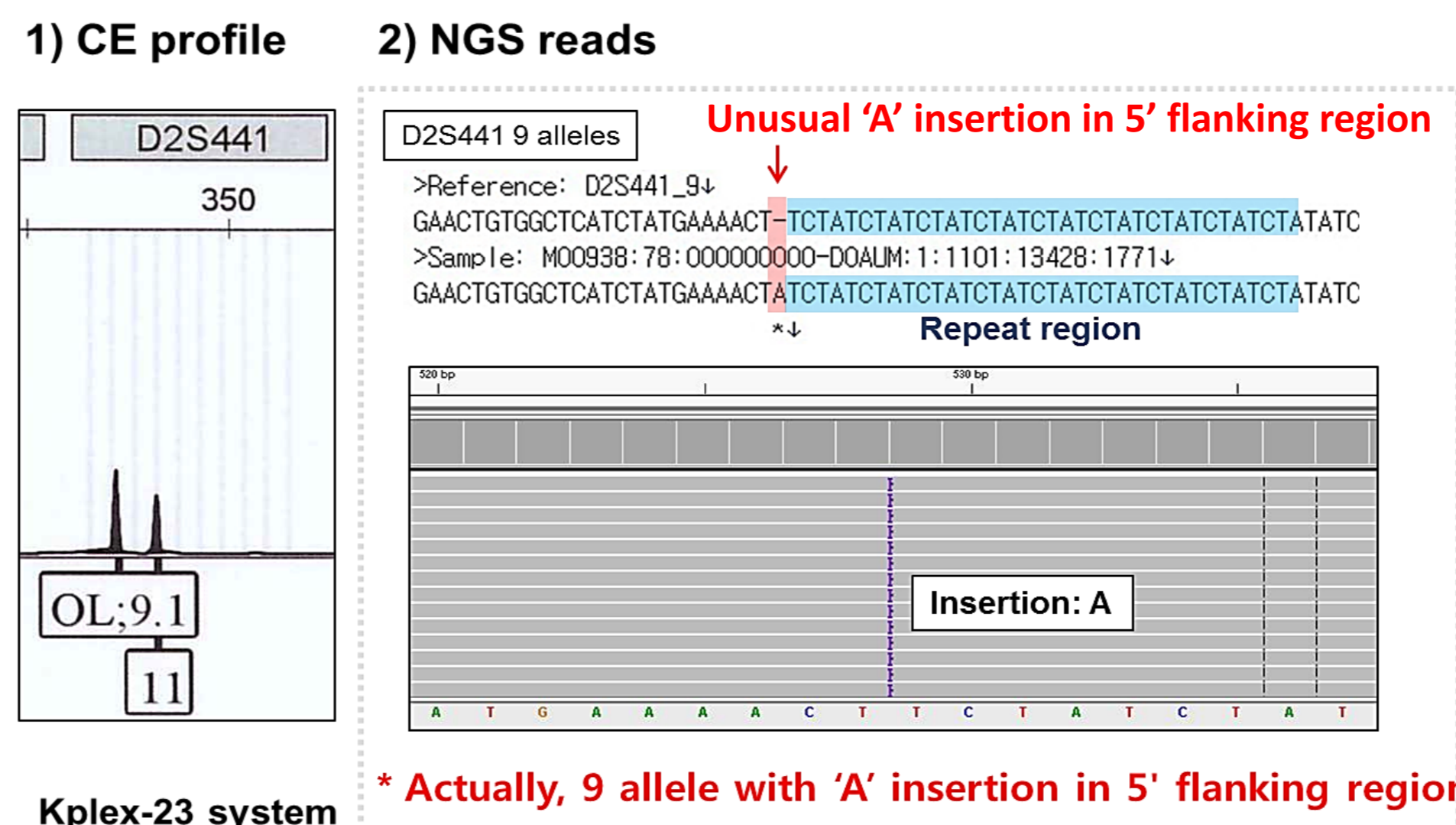
Fig. 3. Read count from NGS data on each marker



The relative reads counts in average coverage by STRait Razor. The minimum and maximum coverage were observed in D5S818 and D6S1043, respectively. The difference was in less than 3 times.

Fig. 4. Comparison of genotype between CE and NGS analysis

STR genotypes of 160 Koreans obtained from NGS were concordant with CE profiles except 14 samples. Fourteen samples were identified that an adenine was inserted at 5' flanking region in 9 allele of D2S441, but shown as 9.1 allele in CE profile.



\* Actually, 9 allele with 'A' insertion in 5' flanking region

2) D2S441

allele	Sub-allele	Structure	Frequency
9	9	[TCTA] <sub>9</sub>	0.041
10	10a	[TCTA] <sub>10</sub>	0.069
	10b	[TCTA] <sub>8</sub> TCTG TCTA	0.147
11	11a	[TCTA] <sub>11</sub>	0.391
	11b	[TCTA] <sub>9</sub> TCTG TCTA	0.009
11.3	11.3	[TCTA] <sub>4</sub> TCA [TCTA] <sub>7</sub>	0.016
12	12	[TCTA] <sub>12</sub>	0.169
13	13a	[TCTA] <sub>13</sub>	0.009
	13b	[TCTA] <sub>10</sub> TTTA [TCTA] <sub>2</sub>	0.013
14	14a	[TCTA] <sub>14</sub>	0.003
	14b	[TCTA] <sub>11</sub> TCTG [TCTA] <sub>2</sub>	0.003
14c	14c	[TCTA] <sub>11</sub> TTTA [TCTA] <sub>2</sub>	0.122
	15	[TCTA] <sub>12</sub> TTTA [TCTA] <sub>2</sub>	0.009

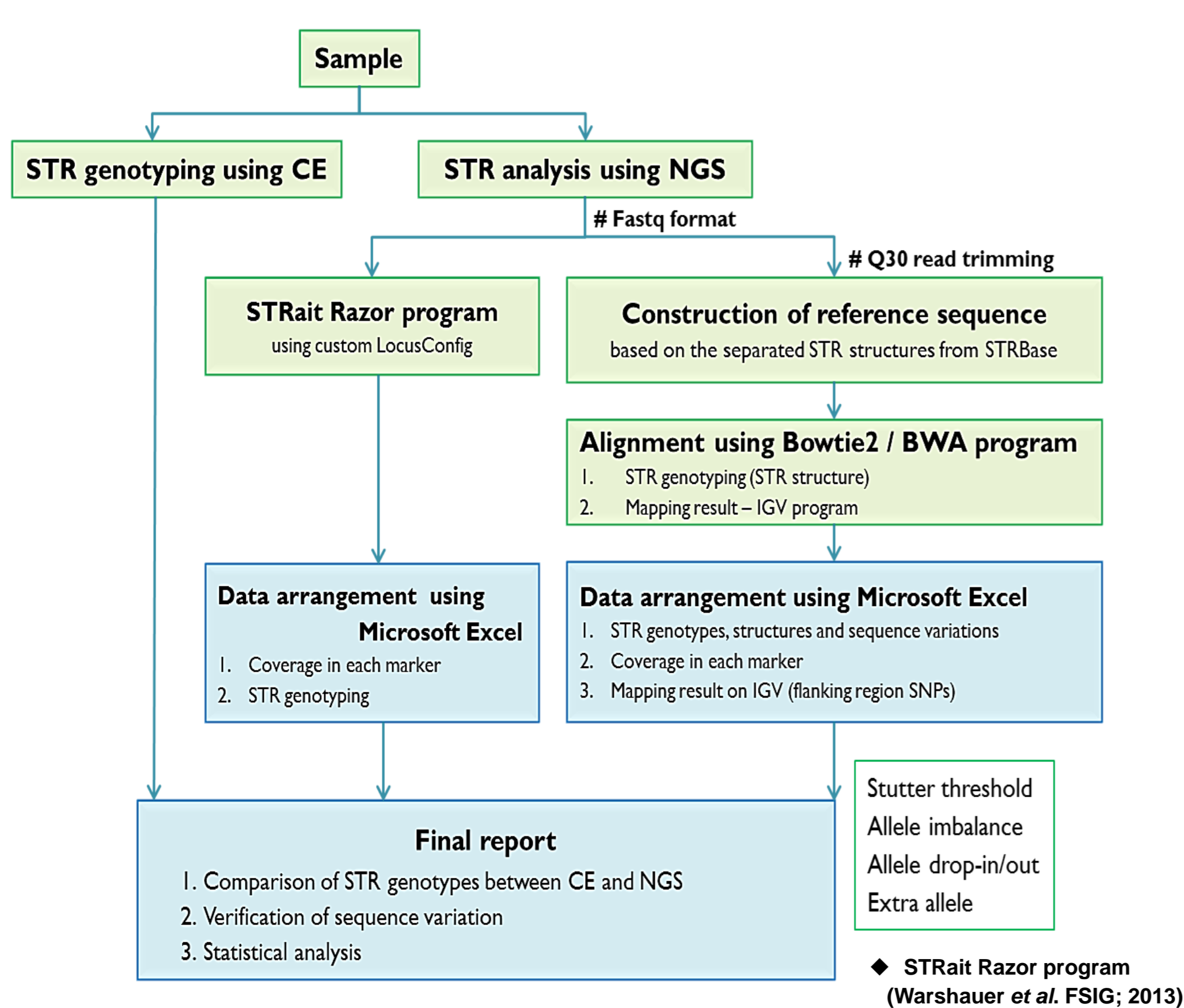
The sequence variations were observed in 12 STRs (D1S1656, D2S441, D2S1338, D3S1358, D6S1043, D7S820, D8S1179, D12S391, D13S317, D21S11, CSF1PO and vWA) of 160 Koreans. D12S391 and D2S441 were presented as examples. And the remaining of sequence variations in 10 STRs will be published in a paper.

Table 2. Comparison between CE and NGS analysis

Loci	Comparison of the number of alleles (n=160)				Comparison of the number of genotypes (n=160)			
	CE method	NGS method	Fold change	Loci	CE method	NGS method	Fold change	
D12S391	12	34	+2.83x	D21S11	31	82	+2.65x	
D21S11	13	29	+2.23x	D3S1358	14	31	+2.21x	
D2S1338	13	29	+2.23x	D12S391	31	68	+2.19x	
D8S1179	9	17	+1.89x	D8S1179	31	65	+2.10x	
D3S1358	8	15	+1.88x	D2S1338	43	80	+1.86x	
D2S441	8	13	+1.63x	D2S441	22	32	+1.45x	
vWA	7	11	+1.57x	vWA	22	29	+1.32x	
D1S1656	13	17	+1.31x	D1S1656	40	50	+1.25x	
D13S317	8	10	+1.25x	D13S317	23	25	+1.09x	
D7S820	7	8	+1.14x	CSF1PO	17	18	+1.06x	
CSF1PO	8	9	+1.13x	D7S820	21	22	+1.05x	
D6S1043	12	13	+1.08x	D6S1043	50	51	+1.02x	

The most variable sequences were observed in D12S391, following by D21S11 and D2S1338. The most variable genotypes were observed in D21S11, and following by D3S1358.

Fig. 2. Pipeline of NGS data analysis



Basically, the genotypes were compared with CE method and NGS data. The percentage coverage values were determined by dividing an assigned coverage for each allele by the total coverage of the locus. STR alleles could be determined when 20% of total coverage was used as a threshold.

Table 1. Examples of observed sequence variations

allele	Sub-allele	Structure	Frequency
15	15	[AGAT] <sub>8</sub> [AGAC] <sub>6</sub> AGAT	0.025
	16	[AGAT] <sub>8</sub> [AGAC] <sub>7</sub> AGAT	0.006
16	16a	[AGAT] <sub>8</sub> [AGAC] <sub>7</sub> AGAT	0.006
	16b	[AGAT] <sub>8</sub> [AGAC] <sub>6</sub> AGAT	0.003
17	17	[AGAT] <sub>10</sub> [AGAC] <sub>6</sub> AGAT	0.066
	18a	[AGAT] <sub>9</sub> [AGAC] <sub>8</sub> AGAT	0.003
18	18b	[AGAT] <sub>10</sub> [AGAC] <sub>7</sub> AGAT	0.025
	18c	[AGAT] <sub>11</sub> [AGAC] <sub>6</sub> AGAT	0.266
18d	18d	[AGAT] <sub>12</sub> [AGAC] <sub>5</sub> AGAT	0.006
	19a	[AGAT] <sub>11</sub> [AGAC] <sub>7</sub> AGAT	0.031
19	19b	[AGAT] <sub>11</sub> [AGAC] <sub>8</sub>	0.003
	19c	[AGAT] <sub>12</sub> [AGAC] <sub>6</sub> AGAT	0.228
19d	19d	[AGAT] <sub>13</sub> [AGAC] <sub>5</sub> AGAT	0.003
	19.2	[AGAT] <sub>4</sub> AT [AGAT] <sub>7</sub> [AGAC] <sub>7</sub> AGAT	0.003
20a	20a	[AGAT] <sub>11</sub> [AGAC] <sub>8</sub> AGAT	0.009
	20b	[AGAT] <sub>11</sub> [AGAC] <sub>9</sub>	0.006
20c	20c	[AGAT] <sub>12</sub> [AGAC] <sub>7</sub> AGAT	0.025
	20d	[AGAT] <sub>12</sub> [AGAC] <sub>8</sub>	0.006
20e	20e	[AGAT] <sub>13</sub> [AGAC] <sub>6</sub> AGAT	0.091
	20f	[AGAT] <sub>13</sub> [AGAC] <sub>7</sub>	0.003
21a	21a	[AGAT] <sub>12</sub> [AGAC] <sub>8</sub> AGAT	0.047
	21b	[AGAT] <sub>12</sub> [AGAC] <sub>9</sub>	0.019
21c	21c	[AGAT] <sub>13</sub> [AGAC] <sub>7</sub> AGAT	0.016
	21d	[AGAT] <sub>13</sub> [AGAC] <sub>8</sub>	0.003
21e	21e	[AGAT] <sub>14</sub> [AGAC] <sub>6</sub> AGAT	0.031
	22a	[AGAT] <sub>12</sub> [AGAC] <sub>10</sub>	0.003
22b	22b	[AGAT] <sub>13</sub> [AGAC] <sub>8</sub> AGAT	0.016
	22c	[AGAT] <sub>13</sub> [AGAC] <sub>9</sub>	0.028
22d	22d	[AGAT] <sub>14</sub> [AGAC] <sub>7</sub> AGAT	0.006
	22e	[AGAT] <sub>15</sub> [AGAC] <sub>6</sub> AGAT	0.006
23a	23a	[AGAT] <sub>12</sub> [AGAC] <sub>10</sub> AGAT	0.003
	23b	[AGAT] <sub>12</sub> [AGAC] <sub>11</sub>	0.003
24a	24a	[AGAT] <sub>14</sub> [AGAC] <sub>9</sub> AGAT	0.003
	24b	[AGAT] <sub>15</sub> [AGAC] <sub>8</sub> AGAT	0.003
25	25	[AGAT] <sub>15</sub> [AGAC] <sub>8</sub> AGAT	0.003

Table 3. Statistical analysis of forensic parameters in Koreans

Loci	Method	$H_{obs}$	$H_{exp}$	MP	PD	PIC	PE
D1S1656	CE	0.831	0.834	0.056	0.944	0.812	0.658
	NGS	0.856	0.846	0.048	0.952	0.827	0.707
D2S441	CE	0.813	0.749	0.108	0.893	0.712	0.622
	NGS	0.825	0.794	0.079	0.921	0.751	0.646
D2S1338	CE	0.844	0.865	0.039	0.961	0.847	0.683
	NGS	0.869	0.907	0.020	0.980	0.896	0.732
D3S1358	CE	0.738	0.712	0.136	0.864	0.658	0.489
	NGS	0.794	0.775	0.080	0.920	0.743	0.587
D6S1043	CE	0.869	0.877	0.033	0.967	0.861	0.732
	NGS	0.869	0.877	0.033	0.967	0.861	0.732
D7S820	CE	0.813	0.769	0.097	0.903	0.730	0.622
	NGS	0.813	0.771	0.095	0.905	0.733	0.622
D8S1179	CE	0.806	0.841	0.047	0.953	0.818	0.611
	NGS	0.881	0.906	0.022	0.978	0.895	0.757
D12S391	CE	0.744	0.801	0.072	0.928	0.772	0.499
	NGS	0.825	0.861	0.037	0.963	0.846	0.646
D13S317	CE	0.781	0.799	0.075	0.925	0.766	0.565
	NGS	0.788	0.801	0.073	0.927	0.769	0.576
D21S11	CE	0.781	0.780	0.080	0.920	0.750	0.565
	NGS	0.919	0.922	0.021	0.979	0.893	0.834
CSF1PO	CE	0.794	0.721	0.146	0.854	0.672	0.587
	NGS	0.794	0.723	0.143	0.857	0.675	0.587
vWA	CE	0.806	0.791	0.080	0.920	0.758	0.611
	NGS	0.806	0.799	0.076	0.924	0.769	0.611

The most significant change of PD between CE and NGS was D21S11 with +0.059. Although the sequence variation was observed in D6S1043, the PD value was not changed.

## Conclusion

- We constructed a single tube new multiplex PCR system that is optimized for NGS analysis of 25 forensic markers with small-sized amplicon.
- STR genotyping results obtained from NGS analysis were mostly consistent with those from CE-based analyses for 160 Koreans except for 14 samples with 9.1 allele of D2S441.
- Sequence variations which differentiate alleles with the same length were observed in 12 STR loci, and the most variable sequences were observed in D12S391.
- Therefore, NGS analysis of global STRs using newly developed multiplex PCR system could provide additional genetic information for forensic investigation.

## 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 (No. 1333-304-260, 2014).