



## Faculty of Health and Medical Sciences

# Sequence variation in the short tandem repeat system SE33 discovered by next generation sequencing

**Eszter Rockenbauer, MSc, PhD and Line Møller, MSc**

**Forensic Geneticist**

**Section of Forensic Genetics**

**Department of Forensic Medicine**

**Faculty of Health and Medical Sciences**

**University of Copenhagen**

**Denmark**



## Aims



- Sequence SE33 with an NGS assay
- Use Roche Junior Platform
- Search for sequence variation
- Identify new allele variants
- Find population specific alleles
- Test usefulness of higher allele diversity in paternity cases



## SE33



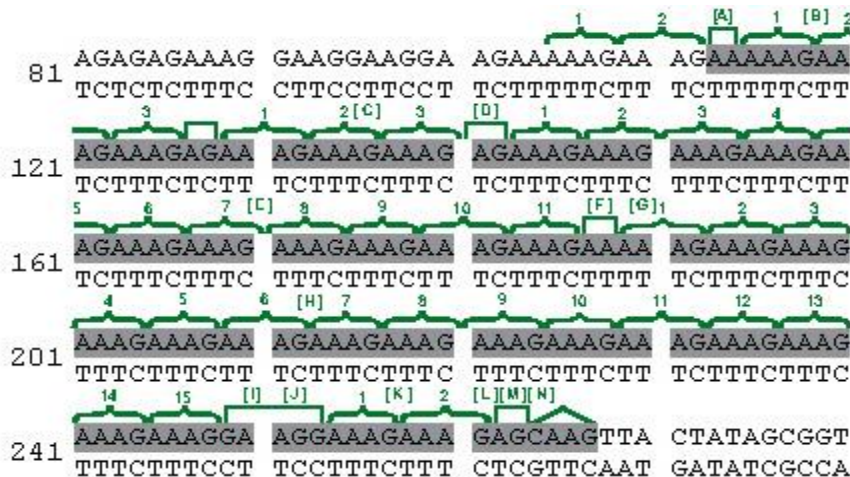
- Core locus in the European Standard Set of STRs
- Included in several commercial STR typing kits
  - E.g. NGM Select, ESI17
- Complex and highly variable sequence
- Consisting of four nucleotide AAAG units interrupted by AA and AG units
- Wide span of allele length variation (3 to 39 four repeat units)





## Challenges with SE33

- No commercially available NGS kits to sequence SE33
- Several polyA stretches difficult to sequence with presently used NGS techniques
- Alleles can be up to 343 bp  $\leftrightarrow$  most NGS machines produce  $\leq 300$  bp reads



## GS Junior 454



- Pyrosequencing technology
- 200 cycles of nucleotide addition (A, T, C and G in fixed order)
- Read length  $\leq 800$  bases
- Accuracy of 99% for reads of  $\leq 400$ bp
- Several successful STR sequencing studies\*

\*Fordyce et al. 2011, Van Neste et al. 2012, Dalsgaard et al. 2014; Rockenbauer et al. 2014; Gelardi et al. 2014; Scheible et al. 2014



# Samples



## Samples

203 samples from Danes(144), Somalis(81) and Greenlanders(58)

- 188 from unrelated individuals
- 5 trios with an unsolved genetic inconsistency in SE33 between parent and child

## Control

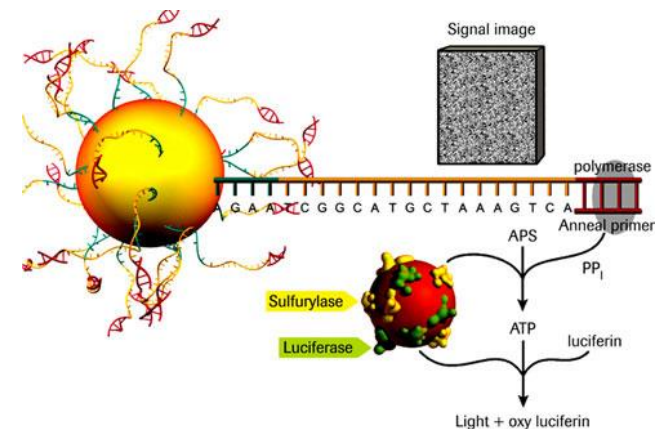
- ISO17025 accredited STR typing (AmpFLSTR<sup>®</sup> NGMSelect<sup>®</sup> PCR Amplification Kit)



# STR sequencing



- Amplicon sequencing with Multiplex IDentifiers (MIDs) in the PCR primer sequence
- Quantification and pooling amplicons into a library
- emPCR and sequencing



## Data analysis



- Generally low sequencing quality
- Standard filters were switched off to get more reads in output file
- "Filters\_off\_amplicons"-pipeline\*:
  - Sorting by MIDs
  - Filtering by flanking sequences
- Both flanking sequence must be present (optional)
- Alignment in BioEdit (Ibis Biosciences)



\*in-house Python-based algorithm





## Results of allele calling



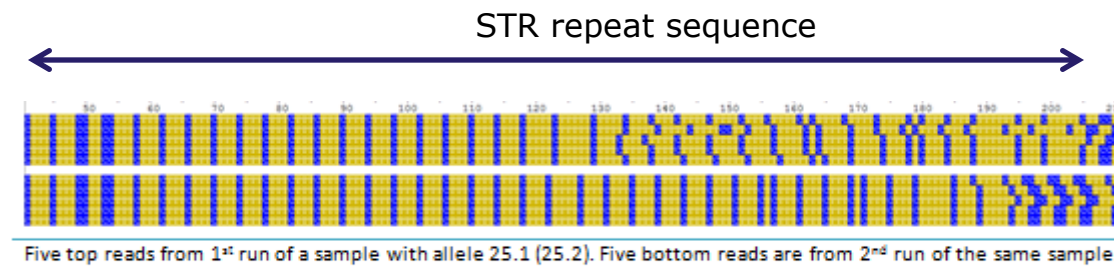
- Total sample coverage 4,643 (before) and 318 (after filtering)
- 10-20% forward 80-90% reverse strand reads
- 57 samples were sequenced twice
- Low-quality samples were reanalysed with only one flanking sequence
- 295 correct allele calls (from 394 expected alleles)
- ~20% called with an insecurity of +/-1bp
- Most correctly called alleles had  $\leq 30$  repeats (~300bp incl. flanks)



## Results of allele calling



- Uncertain calls:
  - Only few reads for long alleles (>30 repeats)
  - Uncertain number of "A's" ("T" in reverse strand read)
- Undetermined alleles:
  - lack of sequencing reads
  - inconsistent reads
  - reads deviating from PCR-CE



## Results of allele calling



- High increase in allele variation
- 27 novel alleles not reported in STRbase
- Raised the power of discrimination by 282% compared to CE
- Identified the mutated allele in 3 out of 5 family trios





## Novel alleles

Table 19 – New Alleles		
1	SE33[13]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>15</sub> G[AAAG] <sub>3</sub> AG
2	SE33[16]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>16</sub> G[AAAG] <sub>3</sub> AG (4bp del in flanking region)
3	SE33[17]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>17</sub> G[AAAG] <sub>3</sub> AG (4bp del in flanking region)
4	SE33[17.3]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> --[AAAG] <sub>15</sub> G[AAAG] <sub>3</sub> G[AAAG] <sub>3</sub> AG
5	SE33[18]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>17</sub> GAAGG[AAAG] <sub>3</sub> AG
6	SE33[21.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AA[AAAG] <sub>13</sub> GAAGG[AAAG] <sub>2</sub> AG
7	SE33[21.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>12</sub> AA[AAAG] <sub>5</sub> GAAGG[AAAG] <sub>2</sub> AG
8	SE33[23]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>23</sub> G[AAAG] <sub>3</sub> AG
9	SE33[24.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>9</sub> <b>A</b> [AAAG] <sub>15</sub> GAAGG[AAAG] <sub>2</sub> AG
10	SE33[25.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>7</sub> <b>A</b> [AAAG] <sub>18</sub> GAAGG[AAAG] <sub>2</sub> AG
11	SE33[25.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>9</sub> AA[AAAG] <sub>17</sub> GAAGG[AAAG] <sub>2</sub> AG
12	SE33[25.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>13</sub> <b>A</b> [AAAG] <sub>12</sub> GAAGG[AAAG] <sub>2</sub> AG
13	SE33[26.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>7</sub> <b>A</b> [AAAG] <sub>19</sub> <b>GAAGG</b> [AAAG] <sub>2</sub> AG
14	SE33[26.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> <b>AG</b> [AAAG] <sub>10</sub> <b>A</b> [AAAG] <sub>17</sub> <b>GG</b> ---[AAAG] <sub>2</sub> AG
15	SE33[26.2]	[AAAG] <sub>2</sub> <b>AG</b> [AAAG] <sub>3</sub> <b>AG</b> [AAAG] <sub>12</sub> <b>A</b> [AAAG] <sub>14</sub> GAAGG[AAAG] <sub>2</sub> AG
16	SE33[27.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>5</sub> <b>A</b> [AAAG] <sub>20</sub> [GAAG] <sub>3</sub> G[AAAG] <sub>2</sub> AG
17	SE33[27.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>5</sub> <b>A</b> [AAAG] <sub>22</sub> GAAGG[AAAG] <sub>2</sub> AG
18	SE33[27.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>7</sub> <b>A</b> [AAAG] <sub>20</sub> GAAGG[AAAG] <sub>2</sub> AG
19	SE33[27.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>9</sub> <b>A</b> [AAAG] <sub>18</sub> GAAGG[AAAG] <sub>2</sub> AG
20	SE33[28.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>9</sub> AA[AAAG] <sub>17</sub> G[AAGG] <sub>3</sub> [AAAG] <sub>2</sub> AG
21	SE33[29.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>9</sub> <b>A</b> [AAAG] <sub>19</sub> G[AAGG] <sub>2</sub> [AAAG] <sub>2</sub> AG
22	SE33[29.2]	[AAAG] <sub>2</sub> <b>AG</b> [AAAG] <sub>3</sub> <b>AG</b> [AAAG] <sub>10</sub> <b>A</b> [AAAG] <sub>19</sub> GAAGG[AAAG] <sub>2</sub> AG
23	SE33[30.2]	<b>A</b> [AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>9</sub> AA[AAAG] <sub>22</sub> GAAGG[AAAG] <sub>2</sub> AG
24	SE33[31.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>11</sub> AA[AAAG] <sub>20</sub> GAAGG[AAAG] <sub>2</sub> (4bp del in flanking region)
25	SE33[30.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>13</sub> <b>A</b> [AAAG] <sub>18</sub> G---[AAAG] <sub>2</sub> AG
26	SE33[31.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>4</sub> <b>A</b> [AAAG] <sub>10</sub> <b>A</b> [AAAG] <sub>17</sub> <b>GAAGG</b> [AAAG] <sub>2</sub> AG
27	SE33[31.2]	[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>8</sub> AG[AAAG] <sub>23</sub> GAAGG[AAAG] <sub>2</sub> AG

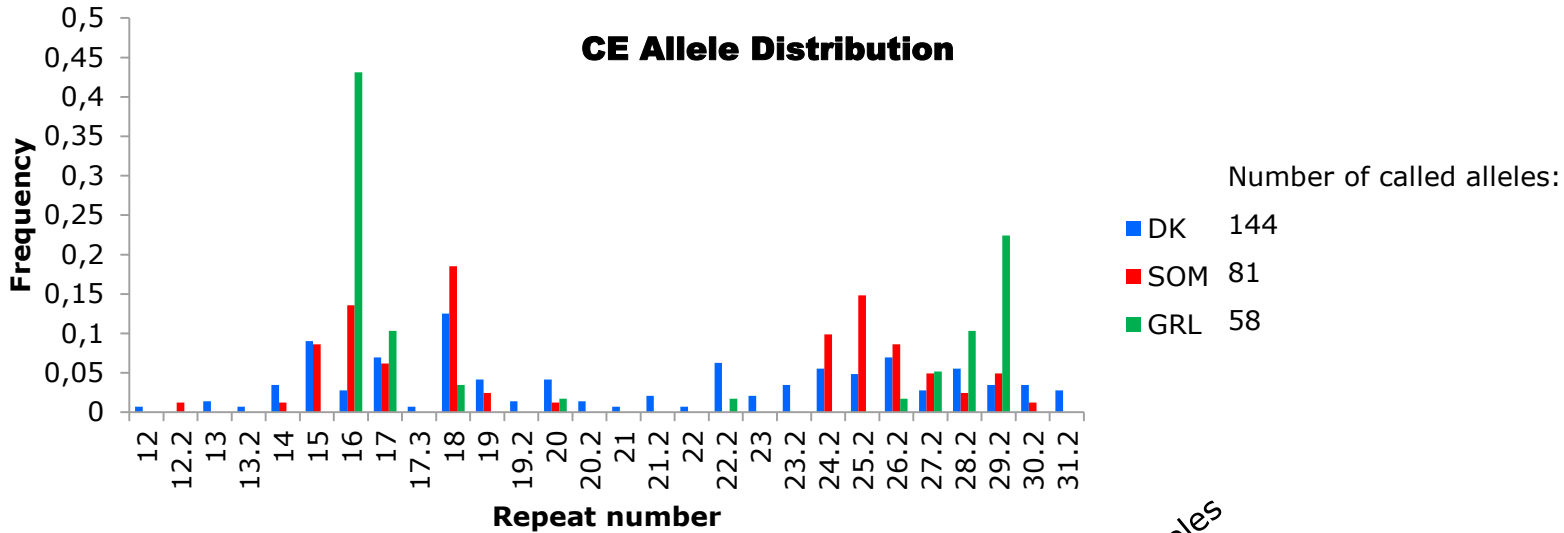
**Bold black** bases are missing in some reads.

**Bold red** bases are appear duplicated in some reads.





# Allele distribution



## PCR-CE / NGS

	Danes		Somali		Greenlanders	
	No. of unique alleles					
Danes	12 43	27 61				
Somali	1 10	14 14	15 28			
Greenlanders	0 4	9 9	8 7	9 15		



# Mutation studies



	Rolle	NGS Allele sequence	CE length
1	M	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>11</sub> AA[AAAG] <sub>10</sub> GAAGG[AAAG] <sub>2</sub>	<u>21.2</u>
1	M	No usable sequence information obtained	(31.2)
1	C	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> <u>AAG</u> [AAAG] <sub>3</sub> AAG[AAAG] <sub>19</sub> G[AAAG] <sub>3</sub>	<u>19</u>
1	C	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AAG[AAAG] <sub>11</sub> A <sub>x</sub> [AAAG] <sub>10</sub> GAAGG[AAAG] <sub>2</sub>	<u>21.2</u>
1	F	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>18</sub> G[AAAG] <sub>3</sub>	18
1	F	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> <u>AAG</u> [AAAG] <sub>3</sub> AAG[AAAG] <sub>20</sub> G[AAAG] <sub>3</sub>	<u>20</u>
2	M	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>17</sub> G[AAAG] <sub>3</sub>	<u>17</u>
2	M	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>18</sub> G[AAAG] <sub>3</sub>	18
2	C	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>17</sub> G[AAAG] <sub>3</sub>	<u>17</u>
2	C	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>15</sub> <u>GAAGG</u> [AAAG] <sub>3</sub>	<u>17</u>
2	F	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>17</sub> <u>GAAGG</u> [AAAG] <sub>3</sub>	<u>18</u>
2	F	No usable sequence information obtained	(20)
5	M	[AAAG] <sub>2</sub> A <sub>x</sub> [AAAG] <sub>3</sub> A <sub>x</sub> G[AAAG] <sub>3</sub> A <sub>x</sub> G[AAAG] <sub>11</sub> A <sub>x</sub> [AAAG] <sub>18</sub> GAAGG[AAAG] <sub>2</sub>	29.2
5	M	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>13</sub> A[AAAG] <sub>18</sub> G[AAAG] <sub>2</sub>	<u>30.2</u>
5	C	[AAAG] <sub>2</sub> A <sub>x</sub> [AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>13</sub> A <sub>x</sub> [AAAG] <sub>15</sub> GAAGG[AAAG] <sub>2</sub>	<u>28.2</u>
5	C	[AAAG] <sub>2</sub> A[AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>13</sub> A[AAAG] <sub>18</sub> G[AAAG] <sub>2</sub>	<u>30.2</u>
5	F	[AAAG] <sub>2</sub> A[AAAG] <sub>2</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>11</sub> A[AAAG] <sub>20</sub> G[AAAG] <sub>2</sub>	30.2
5	F	[AAAG] <sub>2</sub> A <sub>x</sub> [AAAG] <sub>3</sub> AG[AAAG] <sub>3</sub> AG[AAAG] <sub>14</sub> AA[AAAG] <sub>15</sub> GAAGG[AAAG] <sub>2</sub>	<u>29.2</u>

M: mother, F: father, C: child. ( ) Alleles with no sequence information were not sequenced correctly in comparison to CE results and were left out of analysis.



## Conclusions



- SE33 is a promising candidate locus for NGS based genotyping
- Large sequence variation is observed
- Difficulties owing to poly-A stretches and long alleles
- The Roche Junior platform does not perform optimally
- Similar difficulties expected on other available NGS platforms
- Alternative sequencing technique/s and platform/s are needed



# Acknowledgements



Professor and Director of the department:  
Niels Morling

Forensic geneticists:  
Claus Børsting

Bioinformatics:  
Carina G. Jønck

Lab:  
Nadia Jochumsen

**Thank you for your attention**

