



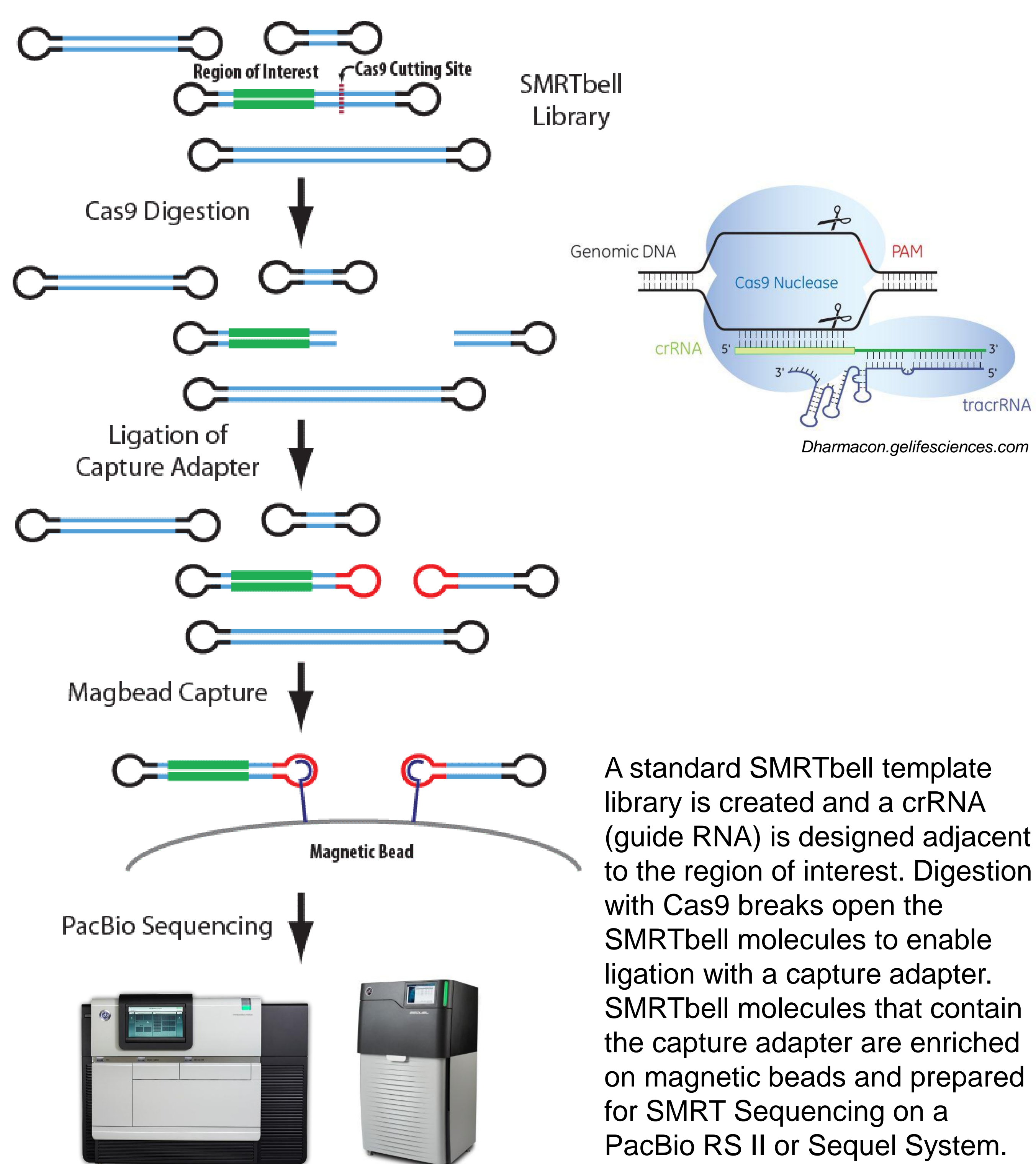
Abstract

Targeted sequencing has proven to be an economical means of obtaining sequence information for one or more defined regions of a larger genome. However, most target enrichment methods rely upon some form of amplification. Amplification removes the epigenetic marks present in native DNA; and some genomic regions, such as those with extreme GC content and repetitive sequences, are recalcitrant to faithful amplification. Yet, a large number of genetic disorders are caused by expansions of repeat sequences. Furthermore, for some disorders methylation status has been shown to be a key factor in the mechanism of disease.

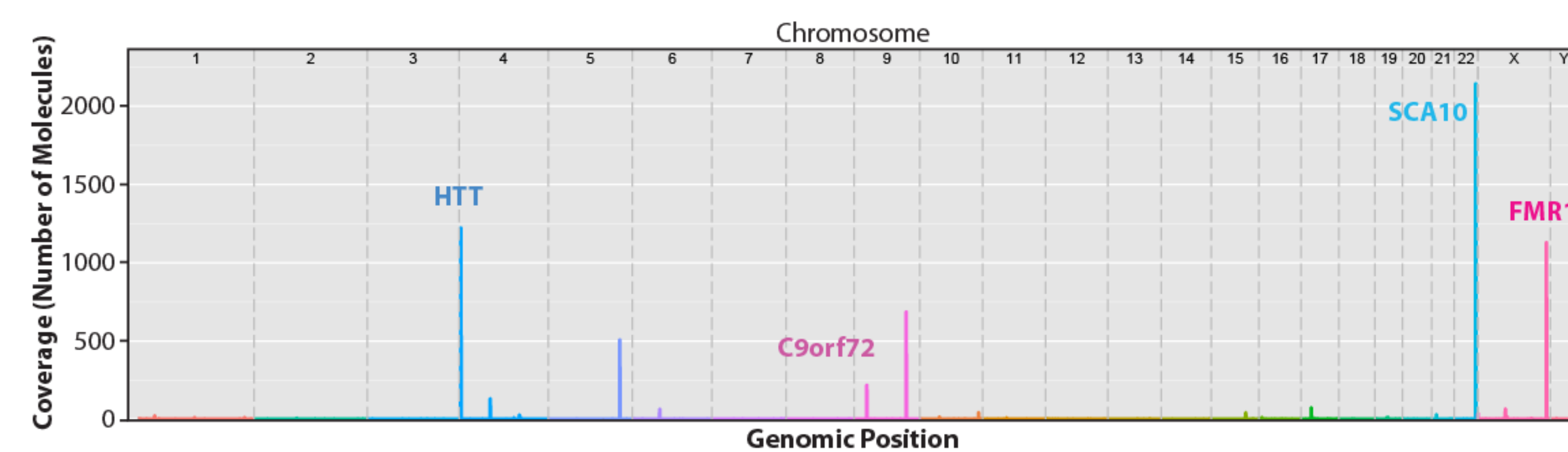
We have developed a novel, amplification-free enrichment technique that employs the CRISPR/Cas9 system for specific targeting of individual human genes. This method, in conjunction with SMRT Sequencing's long reads, high consensus accuracy, and uniform coverage, allows the sequencing of complex genomic regions that cannot be investigated with other technologies. Using human genomic DNA samples and this strategy, we have successfully targeted the loci of a number of repeat expansion disorders (*HTT*, *FMR1*, *SCA10*, *C9orf72*) and disease-associated homonucleotide stretches (*TOMM40*).

With these data, we demonstrate the ability to isolate hundreds of individual on-target molecules and accurately sequence through long repeat stretches, regardless of the extreme GC-content, followed by sequencing on a single PacBio RS II SMRT Cell. The method is compatible with multiplexing of multiple targets and/or multiple samples in a single reaction. Furthermore, because this technique also preserves native DNA molecules for sequencing, it allows for the possibility of direct detection and characterization of epigenetic signatures. We demonstrate detection of 5-mC in the context of repeat expansions.

Method Overview



Targeted Sequencing of 4 Repeat Expansions



| Target Gene | Associated Disease(s) | Chr | crRNA Coordinates | Strand | Target Size | Repeat |
|----------------|--|--------|---------------------|--------|-------------|----------------|
| <i>HTT</i> | Huntington's Disease | Chr 4 | 3075105-3075086 | - | 1125 bp | CAG |
| <i>C9orf72</i> | Familial Frontotemporal Dementia (FTD) and Amyotrophic Lateral Sclerosis (ALS) | Chr 9 | 27572970-27572989 | + | 1261 bp | CCCCGG |
| <i>SCA10</i> | Spinocerebellar Ataxia Type 10 | Chr 22 | 45794847-45794866 | + | 1019 bp | Variable ATTCT |
| <i>FMR1</i> | Fragile X and Fragile X-associated Tremor/Ataxia Syndrome (FXTAS) | Chr X | 147911587-147911606 | + | 1013 bp | CGG |

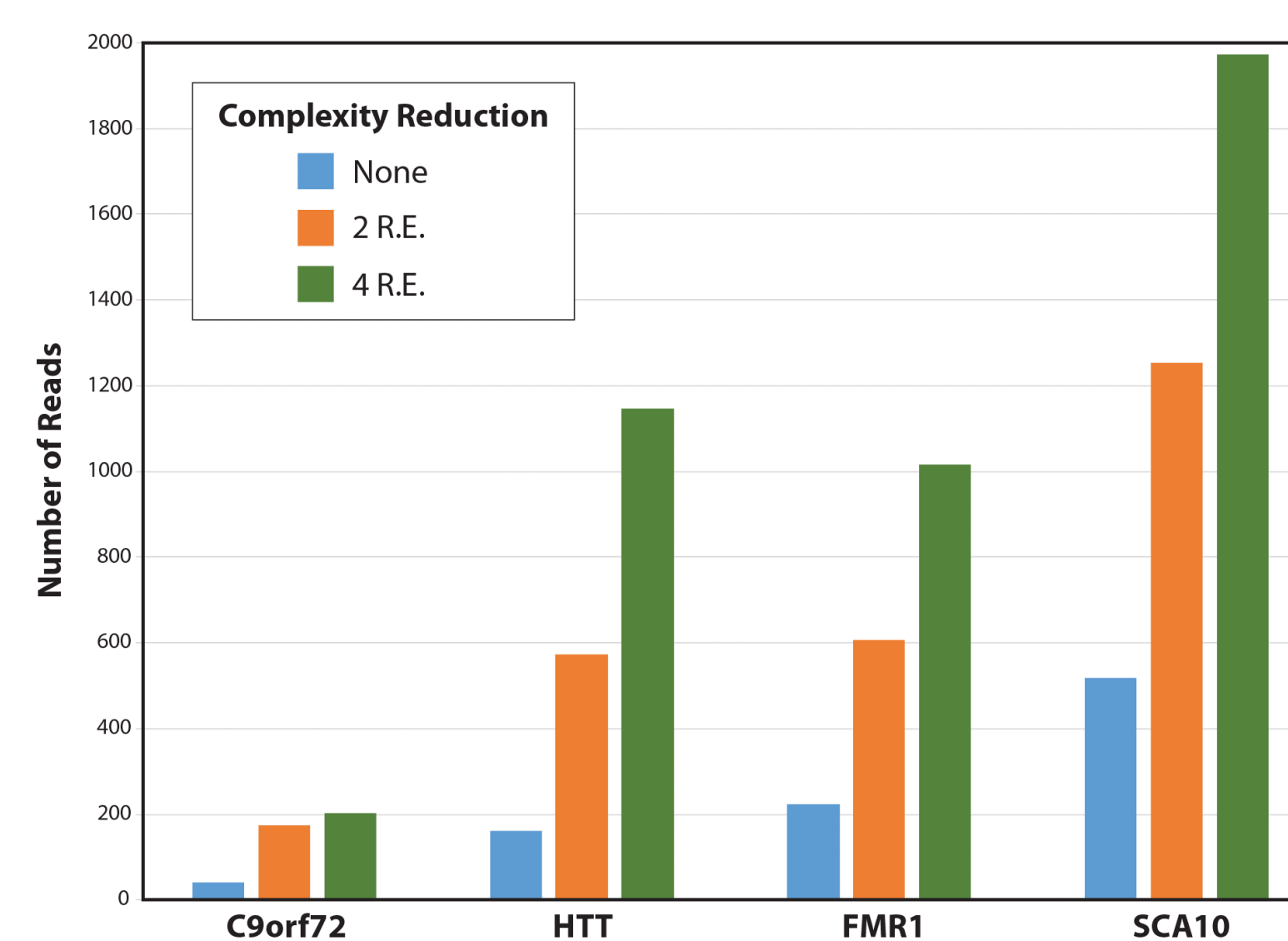
Guide RNAs designed to capture four repeat expansion loci were multiplexed in a single experiment. Molecule coverage across the entire genome is shown above. Off-target signal can be explained by homology of the guide RNA sequence to other regions in the human genome.

Complexity Reduction Improves On-Target Rate

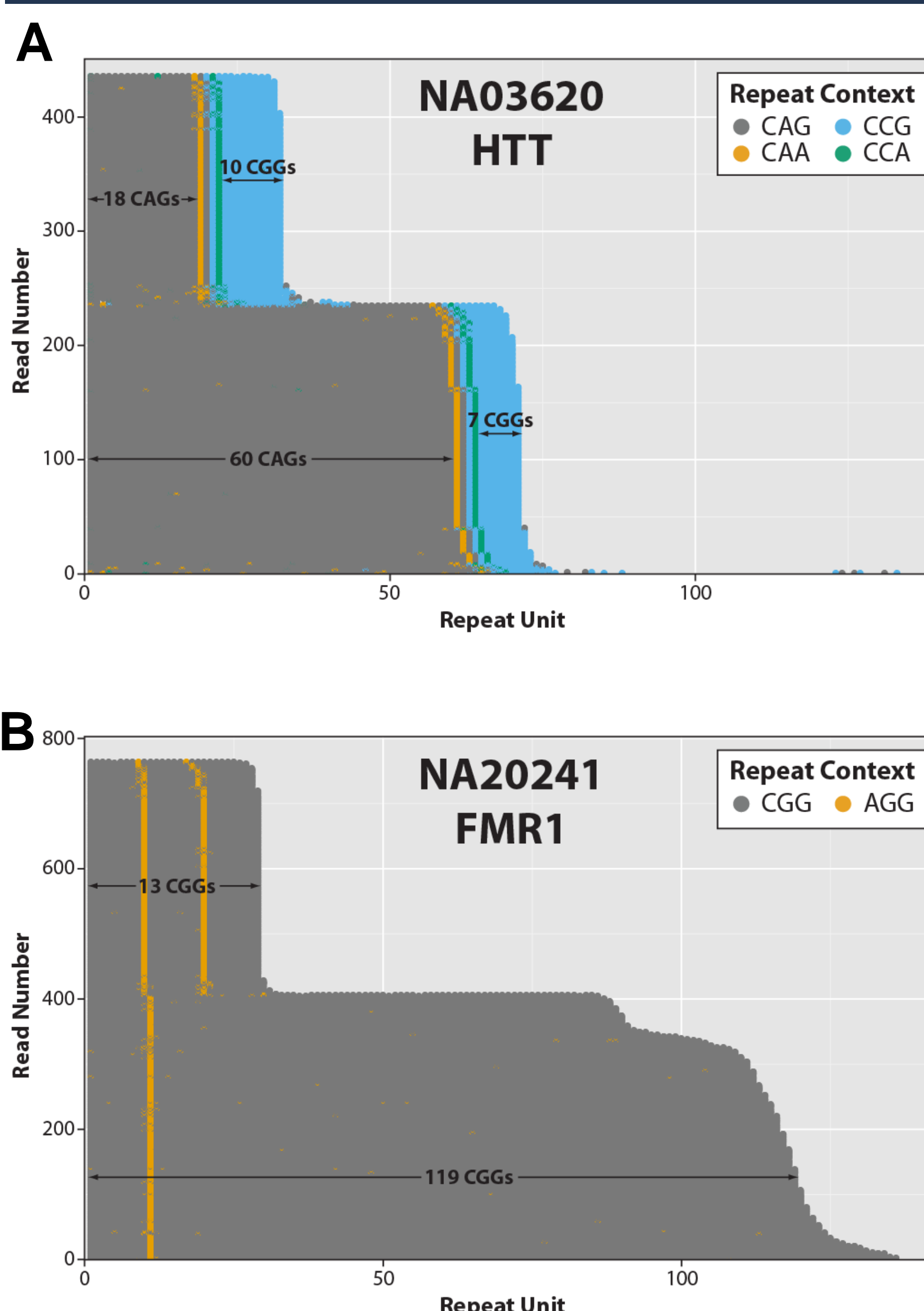
| Complexity Reduction | Input Genomic DNA | Final SMRTbell Yield | % Yield | CCS Reads | On-Target Reads | % Reads On-Target |
|----------------------|-------------------|----------------------|---------|-----------|-----------------|-------------------|
| None | 5.0 µg | 1.9 µg | 37.2% | 44,031 | 945 | 2.15% |
| 2 R.E. | 10.0 µg | 1.5 µg | 15.0% | 51,806 | 2609 | 5.04% |
| 4 R.E. | 20.0 µg | 1.6 µg | 8.0% | 45,676 | 4335 | 9.49% |

0.5 – 1 µg of SMRTbell templates are used as input into the Cas9 reaction

Several restriction enzymes that do not cut within the regions of interest were chosen to remove unwanted SMRTbell templates prior to Cas9 digestion and capture. Inclusion of 2 or 4 restriction enzymes predictably reduces the SMRTbell template yield, but dramatically increases the number of on-target reads and the percentage of reads that come from targeted regions.



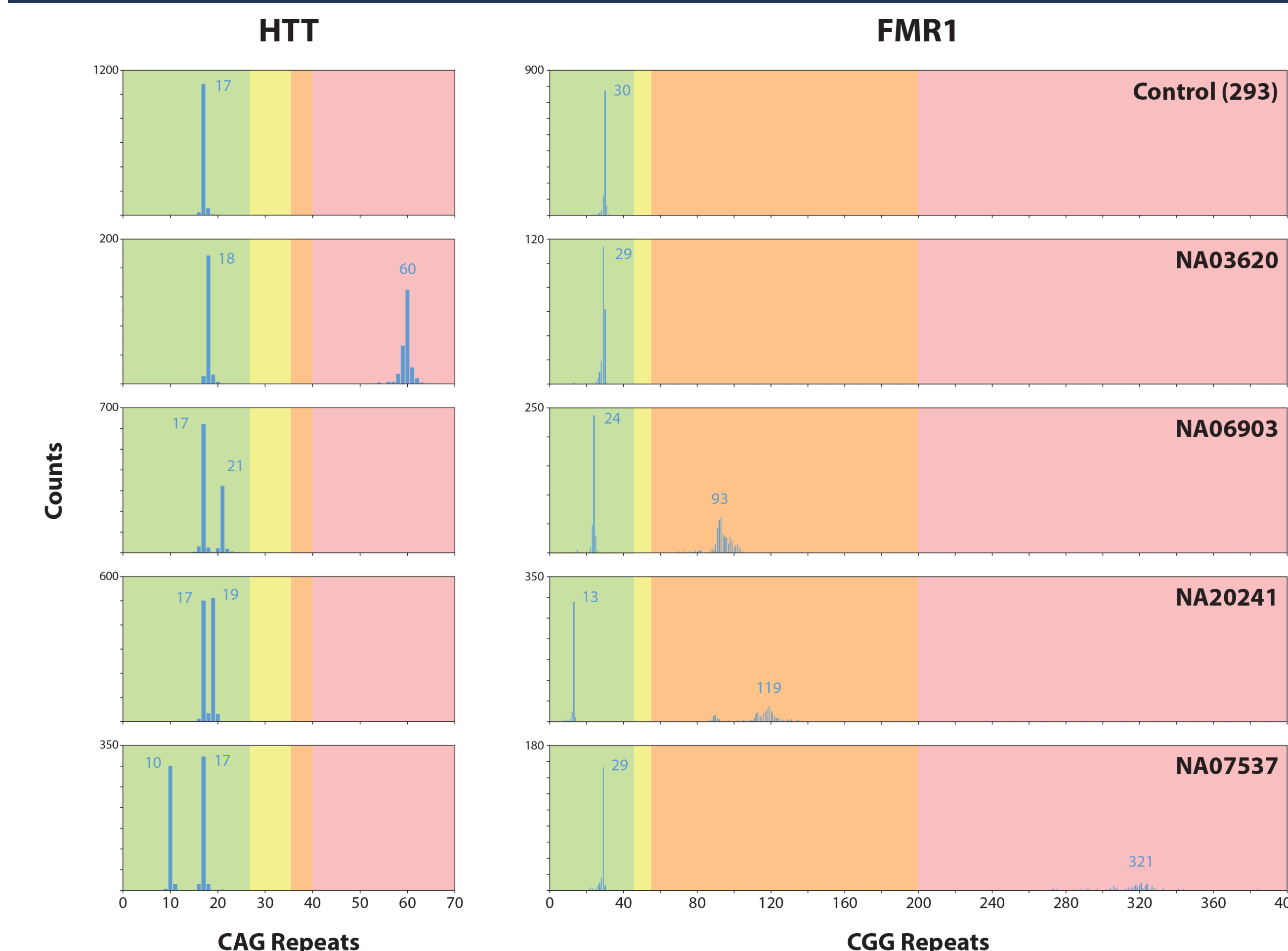
Repeat Structure Visualization



Targeted sequencing of multiple repeat expansion loci was carried out on DNA from two Coriell cell lines from patients with known expansions.

Individual Circular Consensus Sequencing (CCS) reads are trimmed of flanking sequence to include only the relevant repeat region. Trimmed repeat sequences are sorted from shortest to longest. Each individual molecule is represented by a series of colored dots on a horizontal line with each dot representing a single repeat unit, color coded based on the repeat content. (A) *HTT* region in NA03620: Two alleles are visible with varying numbers of CAG and CCG repeats. (B) *FMR1* region in NA20241: Two alleles with varying numbers of CGG repeats and AGG interruptions.

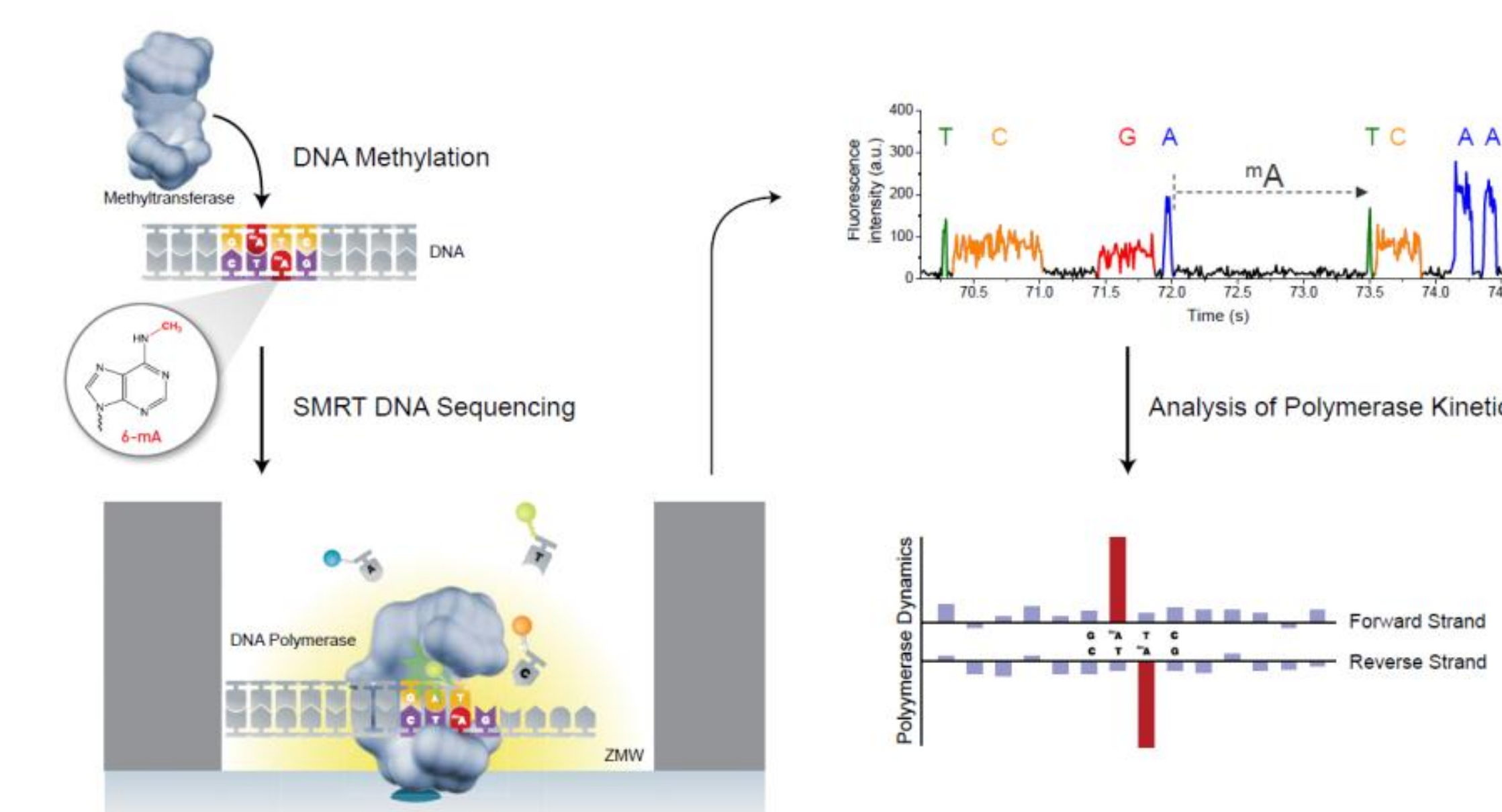
Repeat Count Histograms



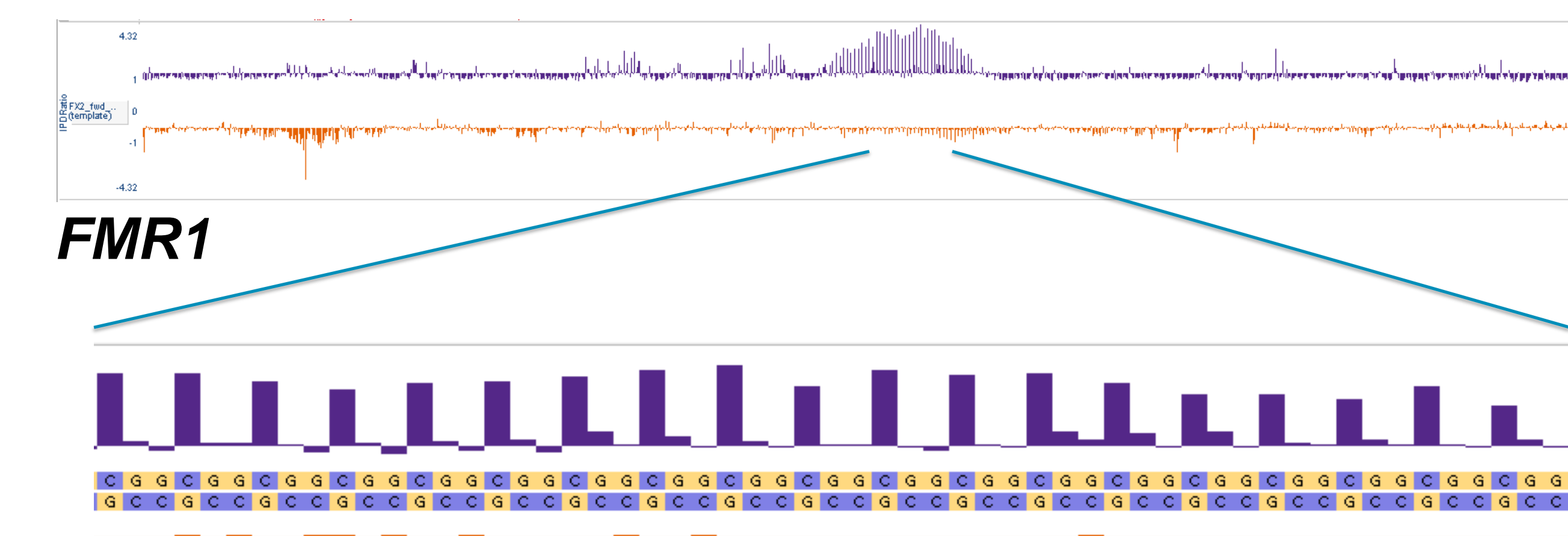
Repeat counts are plotted for the *HTT* (left) and *FMR1* (right) loci across all 5 Coriell samples with count numbers on the y-axis and CAG (*HTT*) or CGG (*FMR1*) repeat numbers on the x-axis. Mode values for each allele are labeled. Shaded background in each plot represents risk ranges for developing disease.

Methylation Detection

Direct Detection of DNA Modifications During SMRT Sequencing



SMRT Sequencing uses kinetic information from each nucleotide to distinguish between modified and native bases.



Kinetic information from a targeted region of the *FMR1* gene shows heavy methylation (5mC) of the CGG repeat.

Summary

Enrich for targeted genomic regions without amplification

- Avoid PCR bias
- Preserve epigenetic modification signals
- Target any genomic region regardless of sequence content

Achieve base-level resolution required to understand the underlying biology of repeat expansion disorder

- Accurately sequence through long repetitive and low-complexity regions
- Count repeats and identify interruption sequences
- Detect mosaicism with single-molecule sequencing