

Proposal for WGBS tracks.

Fractional methylation status is to be published at atleast for all CpGs in the genome. It's open if we want to require publishing this status at other locations, or if we want to leave that optional.

Two tracks are proposed for WGBS data, one indicting the level of methylation, and the other the coverage by C/T at that location (this is essential to have know how much confidence can be had in the methylation call; coverage by other bases is suppressed as that may reveal variants).

The tracks are derived from the following CpGStats file:
(It would be ideal if this file can be shared in the download section of the json hub as because we can visualize only one signal value at a coordinate, we end up with two tracks one for fractional methylation and one for coverage, however, this file contains all information needed for deriving the tracks).

The CpGStats format has three sections:

Bam header section:

This section reports the header of bam underlying the methylation calls as is (except for prefixing each line by #). This ensures that sequence identifiers are present with the methylation calls.

Methylation calling header section:

This is section that includes tool specific comments and the description of data in each column (which must be the last line for this section). At the minimum it consists of:

- Novo5MC 0.1 // software used with version
- java -jar Novo5MC.jar ... // software command line
- chrom position strand #T #C

Data section:

As hinted by the methylation calling header section the columns (which are tab delimited) contain the following information:

- 1: chromosome
- 2: position of C
- 3: strand of C

4: coverageIndicatingUnmethylated [#T] is number of reads showing T at that location [assumption: all unmethylated Cs converted to T]

5: coverageIndicatingMethylated [#C] is number of reads showing C at that location [assumption: no methylated Cs converted to T]

There is one line for every CpG. The line is ignored if $\#C + \#T = 0$ even if there are reads showing alternate bases here (this indicates a variant that is best to not publish).

Based on this file two bedgraph files can be derived: fractional methylation calls and coverage.

The coverage bedgraph contains the sum of last two columns of CpGStats file at each location.

The fractional methylation call reports $10 * \#C / (\#C + \#T)$. Only as many significant figures are reported as number of digits in $\#C + \#T$ (since if $\#C + \#T = 100$, then we can have no more precision than 1/100 in our call).

Both bedgraphs should have header identical to CpGStats file; except for the column description which is now adjusted to:

```
#0-based half-open positioning
#chrom start end fractional_methylation
```

(note from <https://genome.ucsc.edu/goldenpath/help/bedgraph.html> , coordinates are zero-based, half-open: <https://genome.ucsc.edu/FAQ/FAQtracks.html#tracks1>)

The bigwigs from these bedgraphs are proposed as required tracks for WGBS experiments.