

# Single Molecule, Real-Time Sequencing for Base Modification Detection in Eukaryotic Organisms: *Coprinopsis cinerea*



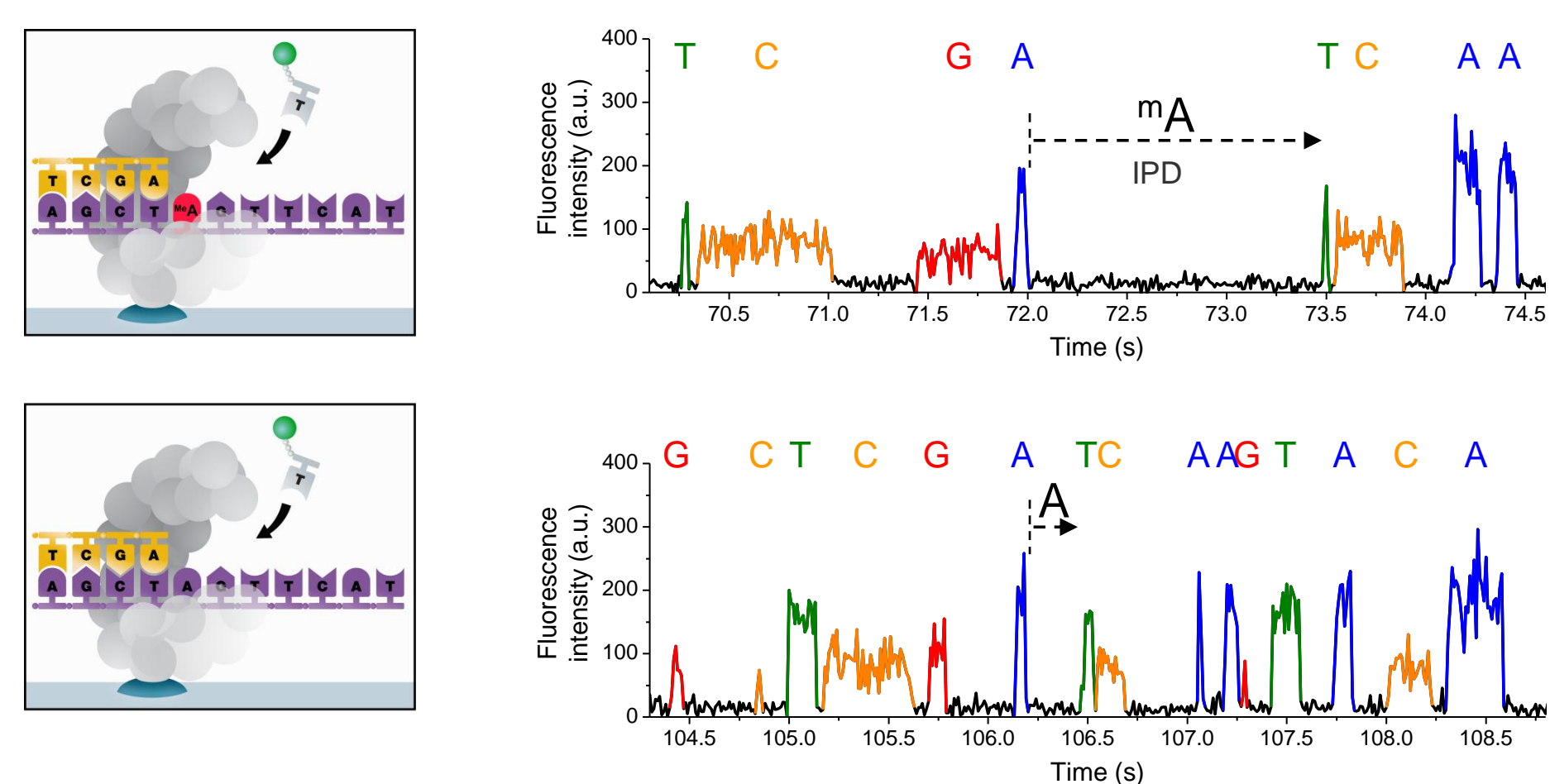
Khai Luong<sup>1</sup>, Tyson A. Clark<sup>1</sup>, Matthew Boitano<sup>1</sup>, Yi Song<sup>1</sup>, Stephen W. Turner<sup>1</sup>, Jonas Korlach<sup>1</sup>  
 Lukas Chavez<sup>2</sup>, Patricia J Pukkila<sup>3</sup>, Yun Huang<sup>2</sup>, Virginia K. Hench<sup>3</sup>, Willaim Pastor<sup>2</sup>, Lakshminarayan M. Iyer<sup>5</sup>, Suneet Agarwal<sup>4</sup>, L. Aravind Iyer<sup>5</sup>, Anjana Rao<sup>2</sup>  
<sup>1</sup>Pacific Biosciences, 1380 Willow Road, Menlo Park, CA 94025  
<sup>2</sup>La Jolla Institute for Allergy & Immunology, 9420 Athena Circle, La Jolla, CA 92037  
<sup>3</sup>University of North Carolina at Chapel Hill, Department of Biology, CB#3280, Chapel Hill, NC 27599-3280  
<sup>4</sup>Harvard Stem Cell Institute, Holyoke Center, Suite 727W, 1350 Massachusetts Ave., Cambridge, MA 02138  
<sup>5</sup>NCBI-CBB, 8600 Rockville Pike, MSC 6075, Bethesda, MD 20894-6075

## Introduction

Single Molecule Real-Time (SMRT<sup>®</sup>) DNA sequencing provides a wealth of kinetic information beyond the extraction of the primary DNA sequence, and this kinetic information can provide for the direct detection of modified bases present in genomic DNA. This method has been demonstrated for base modification detection in prokaryotes at base and strand resolutions. In eukaryotes, the common base modifications known to exist are the cytosine variants including methyl, hydroxymethyl, formyl and carboxyl forms. Each of these modifications exhibits different signatures in SMRT kinetic data, allowing for unprecedented possibilities to differentiate between them in direct sequencing data. We present early results of directly sequencing different base modifications in eukaryotic genomic DNA using this method.

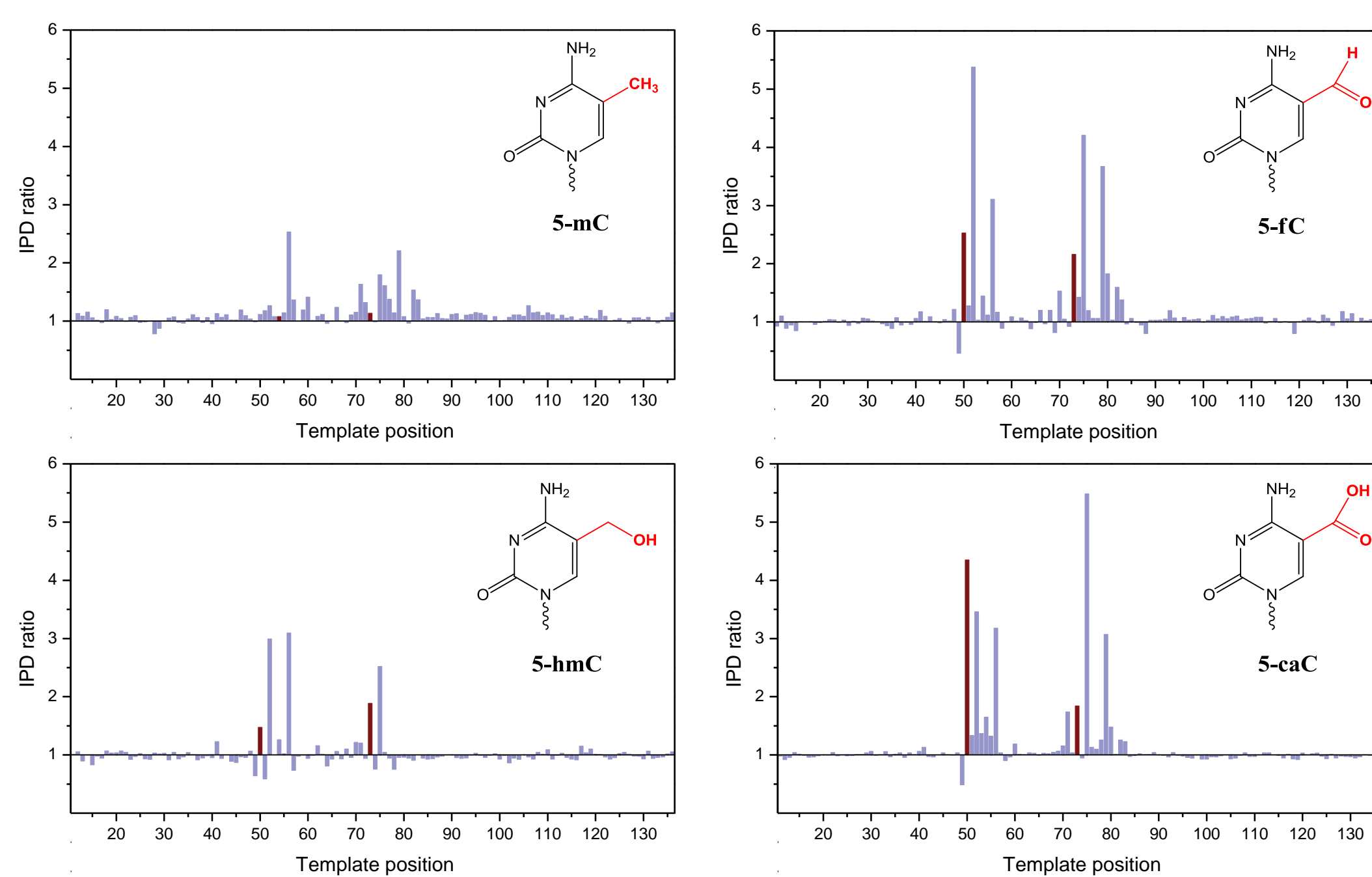
## Methods

### Base Modification Detection by SMRT<sup>®</sup> Sequencing



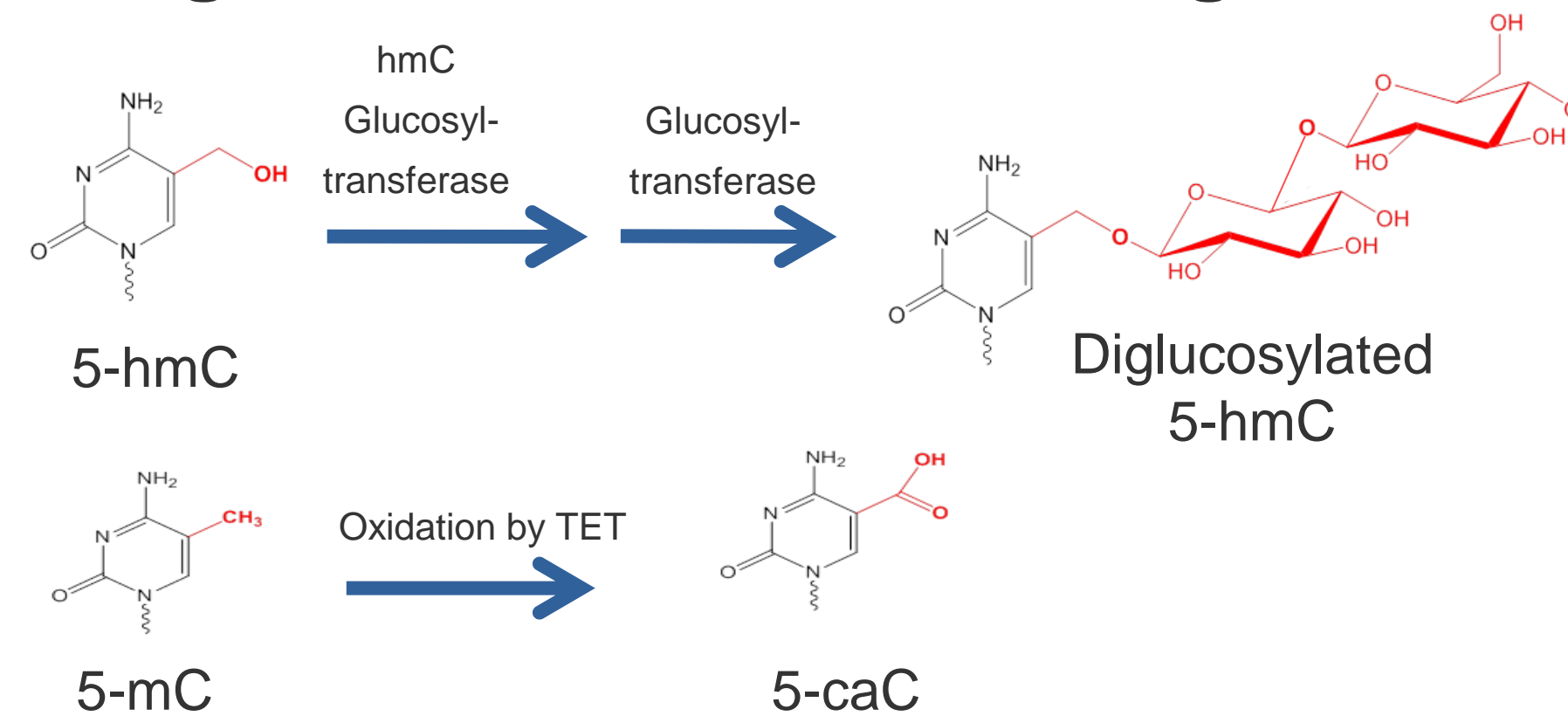
- Inter-pulse duration (IPD) is the time between the previous and current base incorporations
- IPD for a base complement to a modified base is, on average, longer than to a canonical base
- At every position, compare the observed IPDs to the expected IPD distribution

### SMRT<sup>®</sup> Sequencing of the Four Forms of Cytosine



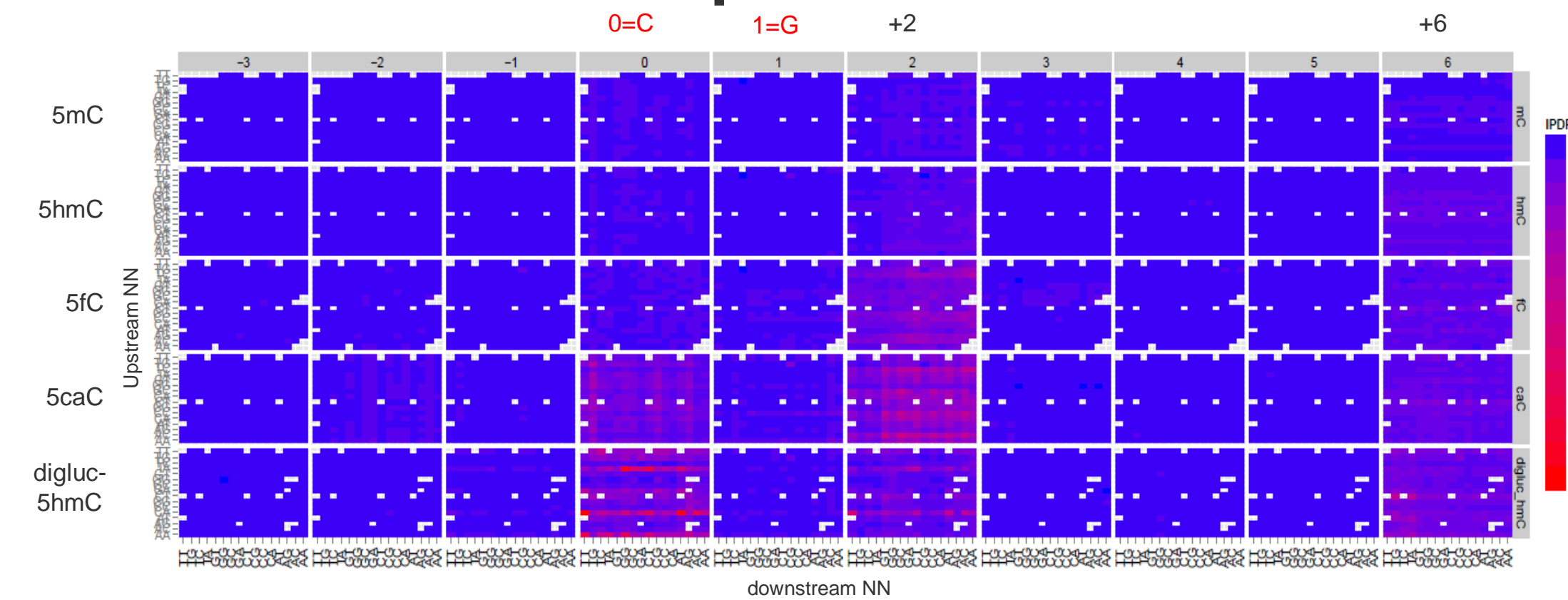
- IPD ratio kinetograms show reproducible footprint signatures for different modifications
- Salient Feature #1: Variants of C have three strongest peaks at position 0, +2, and +6 in the 5' direction

## Signal Enhancement Strategies



- Chemical or enzymatic treatment of the DNA sample can be used to increase signal intensity

## Heat Map of IPD Ratios



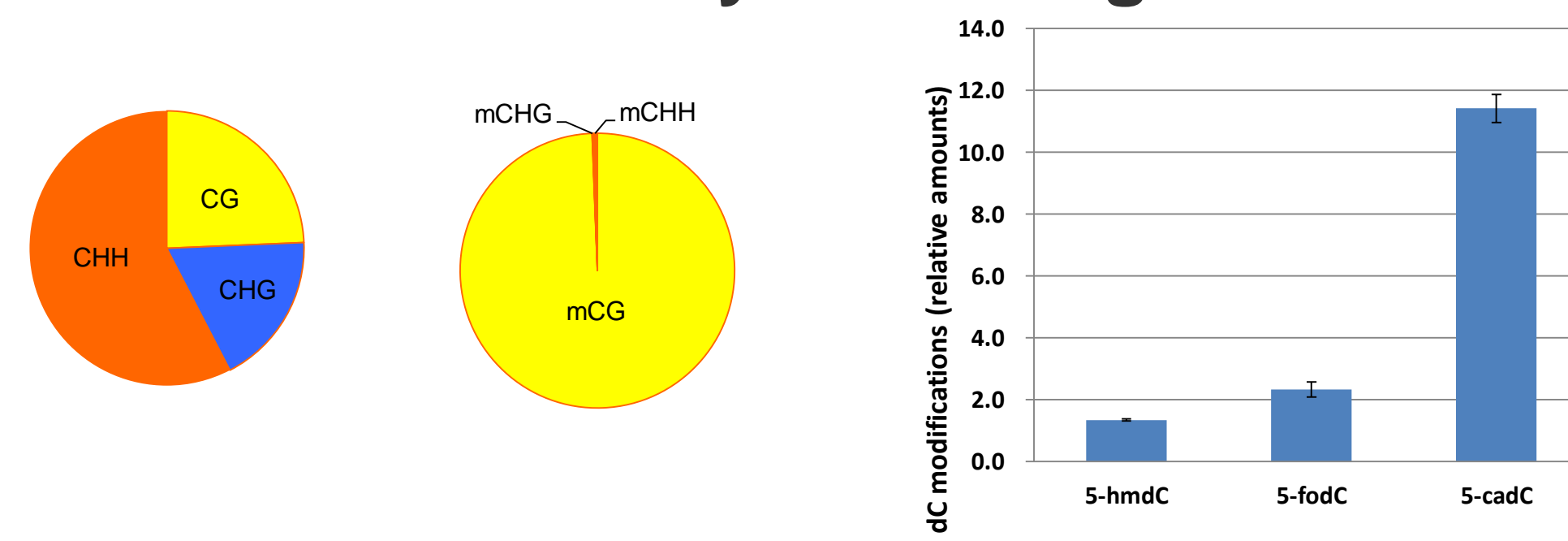
- Using template sequence context NN\*CGNN, the heat map shows sequence-dependent variability in the kinetic signature
- Salient Feature #2: IPD ratio intensity increases from 5-mC to higher oxidative states the of methyl group in nearly all sequence contexts

## Model Organism

- *Coprinus cinereus* (*C. cinerea okayama* 7#130): a multicellular basidiomycete fungus with a typical mushroom form that undergoes a complete sexual cycle.
- Small genome ~36 Mb: 13 chromosomes



## Variants of cytosine in genome



## Results

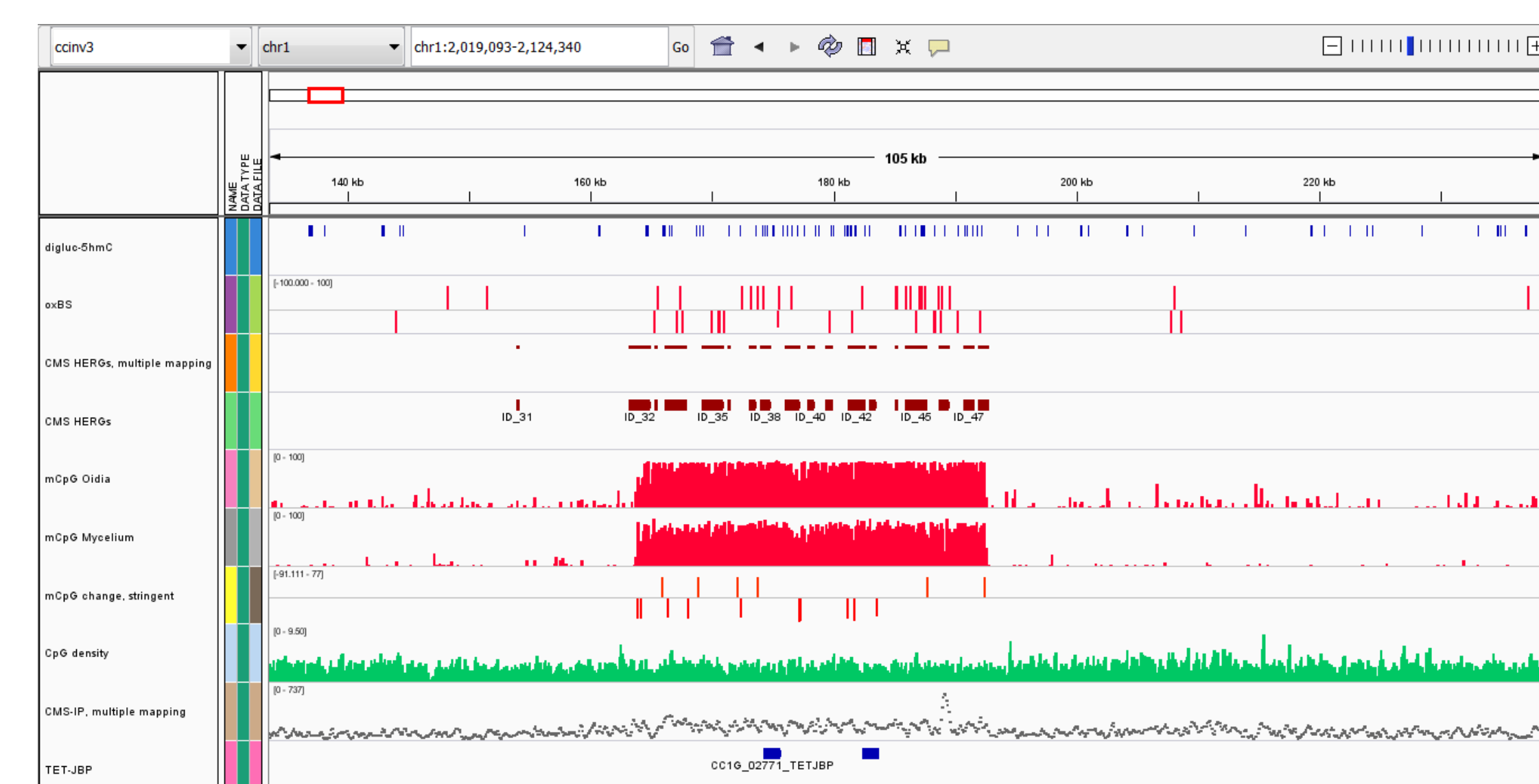
- Used SMRT sequencing to sequence complete *C. cinerea* genome to ~100x coverage using both 800 bp and 8 kb libraries
  - Native genome
  - Diglucosylation (enhance 5-hmC)

## Correlation with other Methods



- Global comparison shows that SMRT detection of diglucosylation agrees with other methods, in particular in CpG-rich regions.

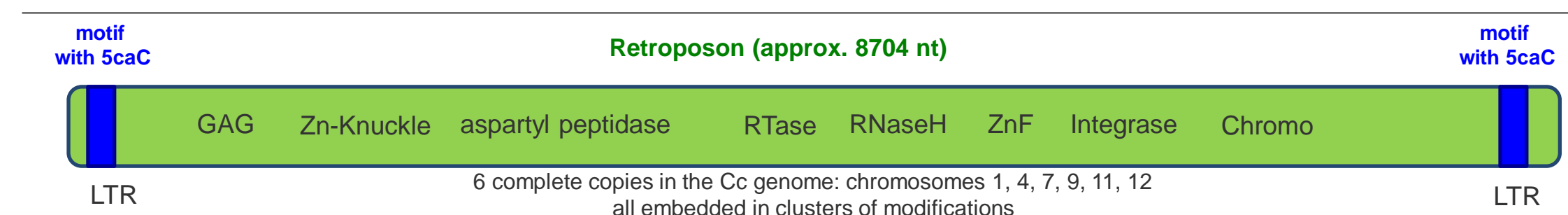
## Correlation with Gene Annotations



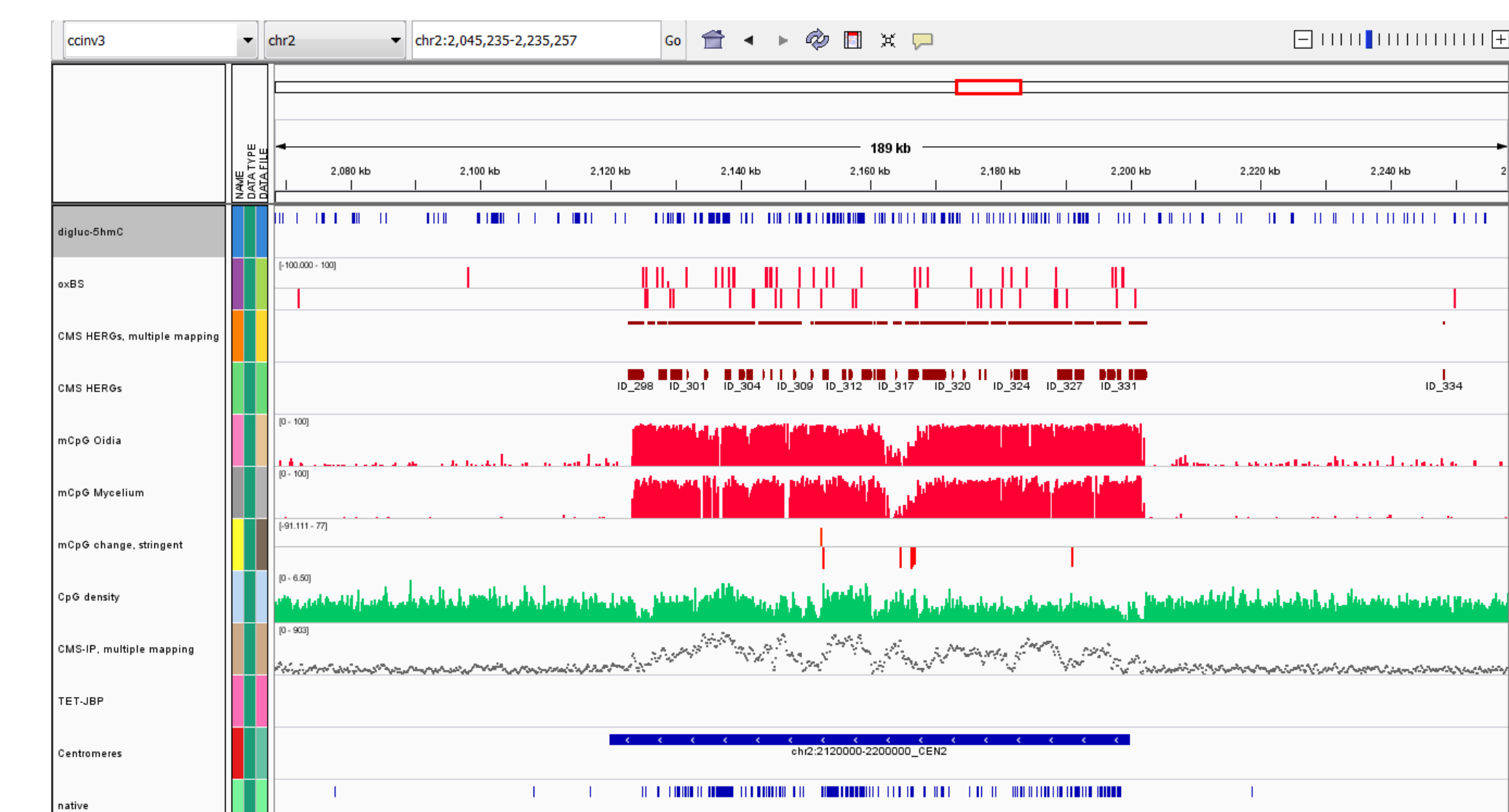
- Modification rich regions overlap with TET/JBP transposons
- Discovered unique motif of ~200 bp that is methylated across several chromosomes and overlaps with novel retroposons

5'...CACAGTTACTGCGGAGCGCAGCAGAGATAAATTAGAGAA...-3

- Comparison of base modification analysis of native vs. diglucosylated samples reveals prominence of caC in LTR of novel retroposons (6 complete copies in chr. 1,4,7,9,11, and 12)



- Pattern of oxidized 5-mC enrichment is predictive for centromeres



## Conclusions

- SMRT Sequencing can detect the four known variants of C found in eukaryote genomes, 5-mC, 5-hmC, 5-fC, and 5-caC
- Enhancement strategy combined with native genome sequencing can increase differentiation between some variants of C
- There are 40 copies of TET/JBP transposons in *C. cinerea*, contained within large regions of oxidized 5-mC in CpG regions
- Centromeres are also strongly marked with oxidized methylcytosines
- A novel set of retroposons bears caC in a specific motif in its long terminal repeat