

Tackling data analysis challenges in genome informatics using HPC

Michael Mueller
Physiological Genomics & Medicine Group
Clinical Genome Informatics Facility
Imperial College London



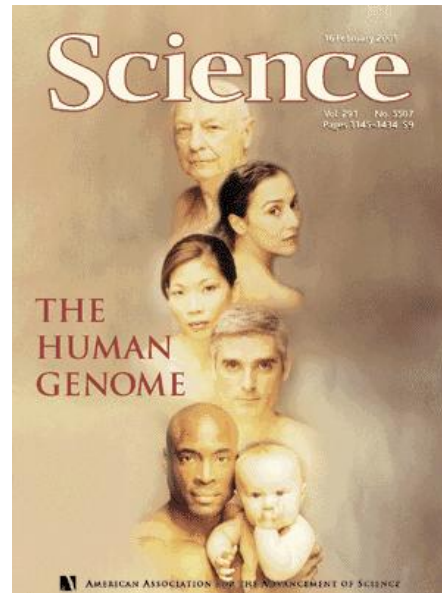
Overview

- Genome sequencing: *why & how?*
- Informatics challenges in genome data analysis
- From genome to epigenome:
Genetic control of inter-individual variation in DNA methylation

The Human Genome

The Human Genome Project 1990 - 2003

determine the sequence of chemical bases (nucleotides) which make up DNA



articles

Finishing the euchromatic sequence of the human genome

International Human Genome Sequencing Consortium*

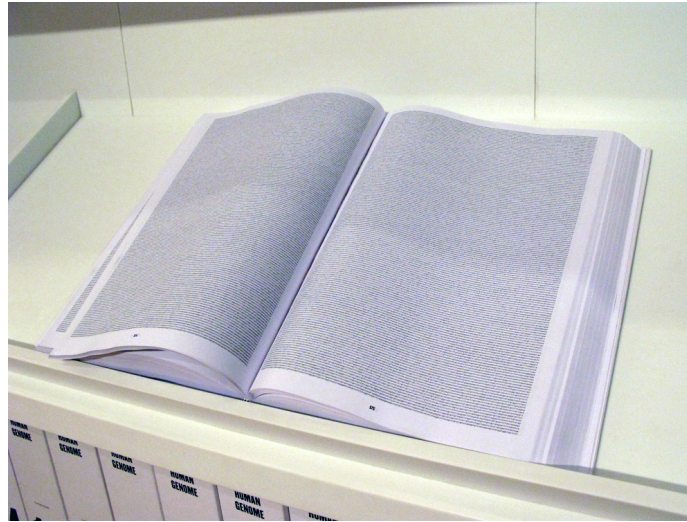
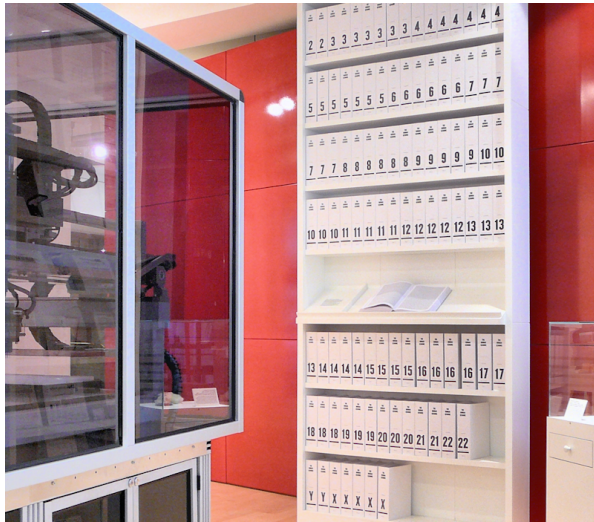
* A list of authors and their affiliations appears in the Supplementary Information

2001

2003

The Human Genome

3.1 billion nucleotides (A,C,T,G)



The Human Genome

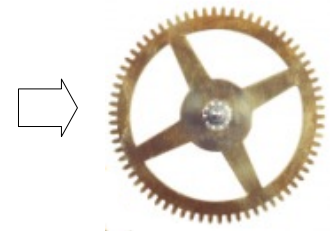
CATGATGAATTCATTTAACCCCTTGAGTTTCCATTTCTTCACCTGCTCCGTGGGGCACT
AACGCCTCCCTCAGAGGCTTCTGGTGAGAATCAGTGTTTCCCTGCCCCCGCCCCGCC
TCCATGCCCCCTTCTCCACGTTCTCACTGTGCTAGGTGCTCTTCTCTGTCTTTCTCTTC
CACCAGCCTGTGGGAAACCTGAGATGAAAGTCGTGTCTTACCCATCTTTGTATTTCCA
GCATCTGAAACTGGGCAGAGCTTAATAAATATTTTTGCTGGAGAGGTTGATGATCTTAC
AAAGTCCCATTGAAAGGTGGCTCTCTGTAAAGCAAAGTTACAATGAGATTGTGATGA
ACATTGTCCCTTGTGGCTTTTCACTTAGTCCCCTCCCTTACCTGAAGAGCAAATTTTC
CTCAAAGTACACAGCAAACAAATGACCCACTGGTGACACTGCTGCCTTTAGACCCCTG
CTGGAAAAGTCTCCACATTTATTAACATTTCCAAAGTAAATTTATCAGGTAGCATTC
ATCAGGTAACATTTGTTGCACATTCATGACTTTTCTACTGTCCACAAAGGCATATGTC
CTTATCATATGCAGACTCCTTGGCCACACTGGATTCCCTCCTTCCCTCCTCGACATGGA
AGAGATGGCATCTTAGGGTCTCTTGTGTTCTTCTGAAGAGGCCTGTGCGGCAGGAAA
AGGCTGCAGCTGCCTTCTGGGAGAAGGAGGAGATGAGTGTATCCTGAACACCTATTA
TGTGCTAGGGGCTATTGTAGATACATGACATTATCATGCTCATTTTTACAAATGAGGA
AACTGAGGCTCAGAAGACTTAAATATTTGCCAAAGAGTTCATAAATGATAGAGCCAG
CATTAGAGTCCAGGACTGTCTGATTTTCAAGCTAAGCTGTTCCCTCTGCACATCATGT
CCCACCAGTAAGGAAGATCTGGGTCTCAGAGCTGAGCCAAGACCTCCCGGGTCTCTG
GGGTTTTCTGTGTCTTTTCAAGTGGCTACATGGAGGTCTGAGAAAAGCCACAGCACA
AAATTGGGCACTGGAGCAAAGAGGAGTGGTGTGGAGGCCTGGCTAAGTATTGACCAAT
GAGCAGGAGTAGGGGCCAGAGCTTGGAGCCCTCAGGTGATAGGTGACCAGGCTGTTGT
TCCACTTTGAAATGCAGGCCCCAGACTAAGGACGGCAGCGAAGCAGAGCTCCCTCGTT
GGTGCAGAGGTTCCCTAAGCTCCTGGGCATTCGTAAGAACTGAGGTTTGCCTTTATTC
TTCTCATGGCTCATGTTATAGCCATTCACCTCCAGAAAGCCTGGCACGTATGGACTGT
TCAAGGCTGTGCTCCATAGAGTAATGATGCCTCCAGCTATGCGAGGCTTTGGGCCAC
CCTTCCACTGCCCCGAGGGCTAAGGGGAGCCCTTCTTCTGCTCCTGCCTGCTCAC
CAGTGTGCCTTTATTAAGTTCGGTGGAGAAACCTCGGTAGGCAGGAGGCATAAGTCCAG
CCACAGAAACCCTGTGCAGATGAGGCTGGGGATGTAGTGAGTGCTGCAGAAGTGAGTG
ACTCAGACACAGAAGAGCTTCGGGTGACAAGCACTAGGACATAGCATTGGATGGGGGG
GAGGTGGGACAAGGAGAGTACTGCCTGTCTGATGTCTGTCTTCCCTAGCTCCAGCT
CTTACCAGATGGGGATCATTGTGGCTGTGGTCACTGGGATTGCTGTAGCGGCCATTGT
TGCTGCTGTAGTGGCCTTGATCTACTGCAGGAAAAAGCGGATTTTCAAGTTTTGTAGCTC
CTCCCGTCCCTTTGGTTATCAGTTTTCACTTGGCCAGGCCCTAACCCAGACATTG
CCAGAATCCCTCTCTTTGGGCTAGATACACATTCAGATCTAGGCCCGTATTGTATTAT
AGTCATTTCATTCGTTTATCTTGTGAGTGGATGACAAAAAGAGGGGAATTGTTAAAGG
AAAAATTTAAATGGAGACTGGAAAAATTCCTGAGCAAAACAAAACACCTGGCCCTTAGA
AATAGCTTTAACTTTGCTTAAACTACAAACACAAGCAAAACTTACGGGGTCATACTA
CATAAAGCATAAGCAAAACTTAACTTGGATGATTTCTGGTAAATGCTTATGTTAGAA
ATAAGACAACCCAGCCAATCACAAGCAGCCTACTAACATATAATTAGGTGACTAGGG

The Human Genome

Functional annotation of the genome sequence

The Human Genome Project
identify and map genes

```
CATGATGAATTCATTTAACCCCTTGAGTTTCCATTTCTTCACCTGCTCCGTGGGGCACT
AACGCCTCCCTCAGAGGCTTCTGGTGAGAATCAGTGTTTCCCTGCCCGCCCGCCGCCC
TCCATGCCCTTCTCCACGTTCTCACTGTGCTAGGTGCTCTTCTCTGTCTTTCTCTTC
CACCAGCCTGTGGGAAACCTGAGATGAAAGTCGTGTCTTACCCATCTTTGTATTTCCA
GCATCTGAAACTGGGCAGAGCTTAATAAATATTTTTGCTGGAGAGGTTGATGATCTTAC
AAAGCTCCCATTGAAAGGTGGCTCTCTGTAAAGCAAAGTTACAATGAGATTGTGATGA
ACATTGTCCCTTGTGGCTTTTCACTTAGTCCCCTCCCTTACCTGAAGAGCAAATTTTC
CTCAAAAGTACACAGCAAACAATGACCCACTGGTGACACTGCTGCCTTTAGACCCCTG
CTGGAAAAGTCTCCACATTTATTAACATTTCCAAAAGTAAATTTATCAGGTAGCATTC
ATCAGGTAACATTTGTTGCACATTCATGACTTTTCTACTGTCCACAAAAGGCATATGTC
CTTATCATATGCAGACTCCTTGGCCACACTGGATTCTCCTTCCCTCCTCGACATGGA
AGAGATGGCATCTTAGGGTCTCTTGTGTTCTTCTGAAGAGGCCTGTCGGGCAGGAAA
AGGCTGCAGCTGCCTTCTGGGAGAAGGAGGAGATGAGTGTATCCTGAACACCTATTA
TGTGCTAGGGGCTATTGTAGATACATGACATTATCATGCTCATTTTTACAAATGAGGA
AACTGAGGCTCAGAAGACTTAAATATTTGCCAAGAGTTCATAAATGATAGAGCCAG
CATTAGAGTCCAGGACTGTCTGATTTTCAAGCTTAAAGTGTTCCTTCTGCACATCATGT
CCCACCAGTAAGGAAGATCTGGGTCTCAGAGCTGAGCCAAGACCTCCCGGGTCTCTG
GGGTTTTCTGTGTCTTTTCAAGTGGCCTACATGGAGGTCTGAGAAAAGCCACAGCACA
AAATTGGGCACTGGAGCAAAGAGGAGTGGTGTGGAGGCCTGGCTAAGTATTGACCAAT
GAGCAGGAGTAGGGGCCAGAGCTTGGAGCCCTCAGGTGATAGGTGACCAGGCTGTTGT
TCCACTTTGAAATGCAGGCCCCAGACTAAGGACGGCAGCGAAGCAGAGCTCCCTCGTT
GGTGCAGAGGTTCCCTAAGCTCCTGGGCATTCGTAAGAAGTGAAGTTTGCCTTTATTC
TTCTCATGGCTCATGTTATAGCCATTCACCTCCAGAAAAGCCTGGCACGTCATGGACTGT
TCAAGGCTGTGCTCCATAGAGTAATGATGCCTCCAGCTATGCGAGGCTTTGGGCCAC
CCTTCCACTGCCCTGAGGGCTAAGGGGAGCCCTTCTTCTGCTCCTGCTGCTCAC
CAGTGTGCCCTTTATTAGTTTGGTGGAGAAACCTCGGTAGGCAGGAGGCATAAGTCCAG
CCACAGAAACCTGTGCAGATGAGGCTGGGGATGTAGTGAGTGTGCAGAAAGTGAAGTG
ACTCAGACACAGAAGAGCTTGGGTGACAAGCACTAGGACATAGCATTGGATGGGGGG
GAGGTGGGACAAGGAGAGTACTGCCTGTCTGATGTCTGTCTTCCCTAGCTCCAGCT
CTTCACCGATGGGGATCATTGTGGCTGTGGTCACTGGGATTGCTGTAGCGGCCATTGT
TGCTGCTGTAGTGGCCTTGATCTACTGCAGGAAAAAGCGGATTTTCAAGTTTTGAGCTC
CTCCCGTCCCTTTGGTTATCAGTTTCCACTTGGCCCAGGCCCTAACCCAGACATTG
CCAGAATCCCTCTCTTTGGGCTAGATACACATTCAGATCTAGGCCCGTATTGTATTAT
AGTCATTCATTGTTTATCTTGTGAGTGGATGACAAAAAGAGGGGAATTGTTAAAGG
AAAAATTTAAATGGAGACTGGAAAAATTCCTGAGCAAAACAAAACCTGGCCCTTAGA
AATAGCTTTAACTTTGCTTAAACTACAAAACAAGCAAAACTTACGGGGTCATACTA
CATAACAAGCATAAGCAAAACTTAACTTGGATGATTTTCTGGTAAATGCTTATGTTAGAA
ATAAGACAACCCAGCCAATCACAAGCAGCCTACTAACATATAATTAGGTGACTAGGG
```



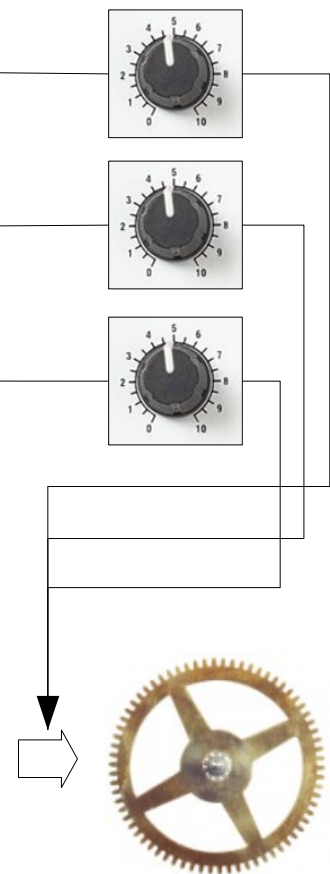
The Human Genome

Functional annotation of the genome sequence

The Human Genome Project
identify and map genes

Encode project (2003)
find functional elements

```
CATGATGAATTCATTTAACCCCTTGAGTTTCCATTTCTTCACCTGCTCCGTGGGGCACT
AACGCCTCCCTCAGAGGCTTCTGGTGAGAATCAGTGTTCCTGCCCCCGCCCCGCC
TCCATGCCCTTCTCCACGTTCTCACTGTGCTAGGTGCTCTTCTCTGTCTTTCTCTTC
CACCAGCCTGTGGGAAACCTGAGATGAAAGTCG|GTCTTACCCAT|CTTTGTATTTCCA
GCATCTGAAACTGGGCAGAGCTTAATAAATATTTTTGCTGGAGAGGTTGATGATCTTAC
AAAGCTCCCATTGAAAGGTGGCTCTCTGTAAAGCAAAGTTACAATGAGATTGTGATGA
ACATTGTCCTTGTGGCTTTTCACTTAGTCCCCTCCCTTCACCTGAAGAGCAAATTTTC
CTCAAAAGTACACAGCAAACAATGACCCACTGGTGACACTGCTGCCTTTAGACCCCTG
CTGGAAAAGAAATCTC|CACATTTATTAA|CAATCCCAAAGTAAATTTATCAGGTAGCATTTC
ATCAGGTAACATTTGTTGCACATTCATGACTTTTCTACTGTCCACAAAGGCATATGTC
CTTATCATATGCAGACTCCTTGGCCACACTGGATTCTCCTTCCCTCCTCGACATGGA
AGAGATGGCATCTTAGGGTCTCTTGTGTTCTTCTGAAGAGGCCTGTCGGGCAGGAAA
AGGCTGCAGCTGCCTTCTGGGAGAAGGAGGAGATGAGTGTATCCTGAACACCTATTA
TGTGCT|AGGGGCTAT|GTAGATACATGACATTATCATGCTCATTTTTACAAATGAGGA
AACTGAGGCTCAGAAGACTTAAATATTTGCCCAAGAGTTTATAAATGATAGAGCCAG
CATTAGAGTCCAGGACTGTCTGATTTTCAAGCCTAAGCTGTTCCCTCTGCACATCATGT
CCCACCAGTAAGGAAGATCTGGGTCTCAGAGCTGAGCCAAGACCTCCCGGGTCTCTG
GGGTTTTCTGTGCTTTTCAGAGTGGCCTACATGGAGGTCTGAGAAAAGGCCACAGCACA
AAATTGGGCACTGGAGCAAAGAGGAGTGGTGTGGAGGCCTGGCTAAGTATTGACCAAT
GAGCAGGAGTAGGGGCCAGAGCTTGGAGCCCTCAGGTGATAGGTGACCAGGCTGTTGT
TCCACTTTGAAATGCAGGCCCCAGACTAAGGACGGCAGCGAAGCAGAGCTCCCTCGTT
GGTGCAGAGGTTCCCTAAGCTCCTGGGCATTCGTAAGAACTGAGGTTTGCCTTTATTC
TTCTCATGGCTCATGTTATAGCCATTCACTCCAGAAAAGCCTGGCAGCTCATGGACTGT
TCAAGGCTGTGCTCCATAGAGTAATGATGCCTCCAGCTATGCGAGGCTTTGGGCCAC
CCTTCCACTG|CCCCTGAGGGCTAAGGGGAGCCCTTCTTCTGCTCCTGCCTGCTCAC
CAGTGTGCCTTTATTAGTT|CGGTGGAGAAACCTCGGTAGGCAGGAGGCATAAGTCCAG
CCACAGAAACCTGTGCAGATGAGGCTGGGGATGTAGTGAGT|GCTGCAGAAGTGAGTG
ACTCAGACACAGAAGAGCTT|CGGGTGACAAGCACTAGGACATAGCATTGGATGGGGGG
GAGGTGGGACAAGGAGAGTACTGCCTGTCTGATGTCTGTCTTCCCTAGCTCCAGCT
CTTACCAGTGGGGATCATTGTGGCTGTGGTCACTGGGATTGCTGTAGCGGCCATTGT
TGCTGCTGTAGTGGCCTTGATCTACTGCAGGAAAAAGCGGATTTTCAAGTTTGTAGCTC
CTCCCGTCCCTTTGGTTATCAGTTTCCACTTGGCCCAGGCCCTAACCCAGACATTG
CCAGAATCCCTCTCTTTGGGCTAGATACACATTAGATCTAGGCCCGTATTGTATTAT
AGTCATTCATTGTTTATCTTGTGAGTGGATGACAAAAAGAGGGGAATTGTTAAAGG
AAAAATTTAAATGGAGACTGGAAAAATTCCTGAGCAAAACAAAACCACTGGCCCTTAGA
AATAGCTTTAACTTTGCTTAAACTACAAAACAAGCAAAACTTACGGGGTCATACTA
CATAAAGCATAAGCAAAACTTAACTTGGATGATTTTCTGGTAAATGCTTATGTTAGAA
ATAAGACAACCCAGCCAATCACAAGCAGCCTACTAACATATAATTAGGTGACTAGGG
```



The Human Genome

Functional annotation of the genome sequence

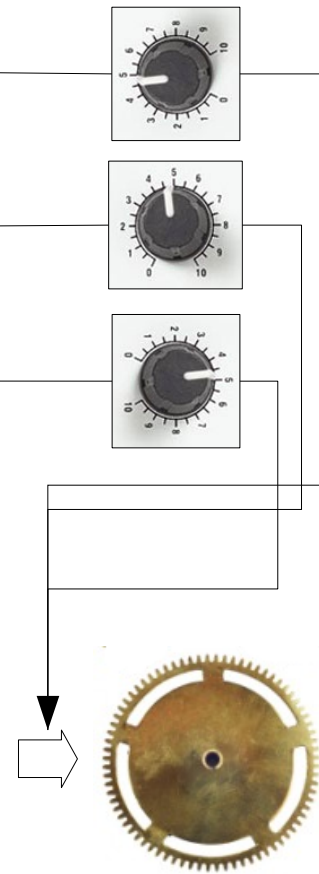
The Human Genome Project identify and map genes

Encode project (2003) find functional elements

1000 genomes project (2008) identify inter-individual sequence variation

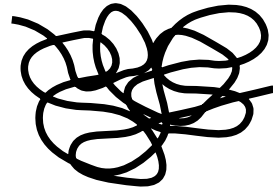


```
CATGATGAATTCATTTAACCCCTTGAGTTTCCATTTCTTCACCTGCTCCGTGGGGCACT
AACGCCTCCCTCAGAGGCTTCTGGTGAGAATCAGTGTTTCCCTGCCCCCGCCCCGCC
TCCATGCCCTTCTCCACGTTCTCACTGTGCTAGGTGCTCTTCTCTGTCTTTCTCTTC
CACCAGCCTGTGGGAAACCTGAGATGAAAGTCG|T|G|C|T|A|C|C|C|A|T|C|T|T|T|G|T|A|T|T|T|C|C|A
GCATCTGAAACTGGGCAGAGCTTAATAAATATTTTTGC|T|G|G|A|G|A|G|T|T|G|A|T|G|A|T|C|T|T|A|C
AAAGCTCCCATTGAAAGGTGGCTCTCTGTAAAGC|T|A|C|A|T|G|A|G|A|T|T|G|T|G|A|T|G|A
ACATTGTCCCTTGTGGCTTTTCACTTAGTCCCCTC|T|A|C|T|G|A|A|G|A|G|C|A|A|A|T|T|T|C
CTCAAAAGTACACAGCAAAACAAATGACCCACTGGTGACACTGCTGCCTTTAGACCCCTG
CTGGAAAAGAAATCTC|C|A|C|A|T|T|T|A|T|T|A|C|A|T|T|T|T|A|T|C|A|G|G|T|A|G|C|A|T|T|C
ATCAGGTAACATTTGTTGCACATTCATGACTTTTCTACTGTCCACAAAGGCATATGTC
CTTATCATATGCAGACTCCTTGGCCACACTGGATTCTCCTTCCCTCCTCGACATGGA
AGAGATGGCATCTTAGGGTCTCTTGTGTTCTTCTGAAGAGGCCTGTGCGGCAGGAAA
AGGCTGCAGCTGCC|T|T|C|T|G|G|G|A|A|G|G|A|G|A|G|A|T|G|A|G|T|G|A|T|C|C|T|G|A|C|A|C|C|T|A|T|T|A
TGTGCT|A|G|G|G|C|T|A|T|T|G|T|A|G|A|T|A|C|A|T|G|A|C|A|T|T|A|T|C|A|T|G|C|A|T|T|T|T|C|A|C|A|A|T|G|A|G|G|A
AACTGAGGC|C|A|G|A|A|G|A|C|T|T|A|A|A|T|A|T|T|T|G|C|C|A|A|G|A|G|T|T|C|A|T|A|A|T|G|A|T|A|G|A|G|C|C|A|G
CATTAG|G|A|C|T|G|A|C|T|G|A|T|T|T|C|A|G|A|C|C|T|A|A|G|C|T|G|T|T|C|C|C|T|C|T|G|C|A|C|A|T|C|A|T|G|T
CCCACC|G|T|A|A|G|A|T|C|T|G|G|G|T|C|A|G|A|G|C|T|G|A|G|C|C|A|A|G|A|C|C|C|C|G|G|G|T|C|C|T|C|T|G
GGGTTTTCTGTG|T|C|T|T|T|C|A|G|A|G|T|G|G|C|T|A|C|A|T|G|G|A|G|G|T|C|T|G|A|G|A|A|A|G|G|C|C|A|G|C|A|C|A
AAATTGGGCACTGGAGCAAAGAGGAGTGGTGTGGAGGCCTGGCTAAGTATTGACCAAT
GAGCAGGAGTAGGGGCCAGAGCTTGGAGCCCTCAGGTGATAGGTGACCAGGCTGTTGT
TCCACTTTGAAATGCAGGCCCCAGACTAAGGACGGCAGCGAAGCAGAGCTCCCTCGTT
GGTGCAGAGGTTCCCTAAGCTCCTGGGCATTGTAAGAAGTGAAGTT|C|G|G|A|T|T|C|A|G|G|T|T|C|A|T|G|G|C|T|C|A|T|G|T|T|A|T|A|G|C|C|A|T|T|C|A|C|T|C|C|A|G|A|A|A|C|C|T|G|G|C|A|G|C
CTGT
TCAAGGCTGTGCTCCATAGATAATGATGCCTCCAGCTATGCGAGGCTTTGGGCCAC
CCTTCCCAC|T|G|C|C|C|T|G|A|G|G|G|C|T|A|A|G|G|G|A|G|C|C|T|T|C|T|T|C|T|G|C|T|C|T|G|C|T|G|C|T|C|A|C
CAGTGTGCC|T|T|A|T|T|A|G|T|T|C|G|G|T|G|G|A|A|A|C|C|T|G|G|T|A|G|G|C|A|G|G|A|G|C|A|A|A|G|T|C|C|A|G
CCACAGAAACC|T|G|T|G|C|A|G|A|T|G|A|G|G|C|T|G|G|G|A|T|G|T|A|G|T|G|A|G|T|G|C|T|G|C|A|G|A|A|G|T|G|A|G|T|G
ACTCAGACACAGAAGAGCTT|C|G|G|G|T|G|C|A|A|G|C|A|C|T|A|G|G|A|C|A|T|G|G|A|T|G|G|G|G|G|G
GAGGTGGGACAAGGAGAGTACTGCCTGTCTGATGTCTGT|G|C|A|G|C|T|G|A|G|C|C|C|A|G|C|T
CTTACCAGATGGGGATCATTGTGGCTGTGGTCACTGGGAT|T|G|C|T|G|T|A|G|C|G|G|C|C|A|T|T|G|T
TGCTGCTGTAGTGGCCTT|G|A|T|C|T|A|C|T|G|C|A|G|G|A|A|A|A|G|C|G|G|A|T|T|C|A|G|G|T|T|G|T|A|G|C|T|C
CTCCCGTCCCTTTGGTTATCAGTTTCCACTTGGCCAGGCCCTAACCCAGACATTG
CCAGAATCCCTCTCTTTGGGCTAGATACACATT|C|A|G|A|T|C|T|A|G|G|C|C|G|T|A|T|T|G|T|A|T|T|A|T
AGTCATT|C|A|T|T|C|G|T|T|A|T|C|T|T|G|C|T|G|A|G|T|G|G|A|T|G|A|C|A|A|A|A|A|G|A|G|G|G|A|A|T|T|G|T|A|A|A|G|G
AAAATTTAAATGGAGACTGGAAAAATTCCTGAGCAAAACAAAACCACCTGGCCCTTAGA
AATAGCTTTAACTTTGCTTAAACTACAAAACAAAGCAAAACTTACGGGGTCATACTA
CATAACAAGCATAAGCAAAACTTAACTTGGATGATTTCTGGTAAATGCTTATGTTAGAA
ATAAGACAACCCAGCCAATCACAAGCAGCCTACTAACATATAATTAGGTGACTAGGG
```

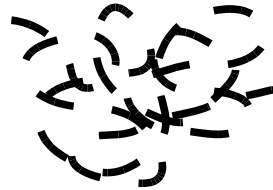


Variant detection by next-generation sequencing

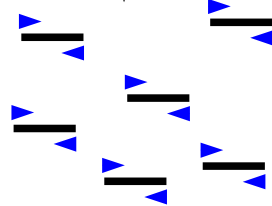
genomic
DNA



fragmented
DNA



paired-end
sequencing



read
mapping

..CTACAT GAGAAT...CACA--A ACTGGA.....TGGCGT GGCTAA.....ACGAGC
GCCTACATGGAGGCTCTGAGAAAGGCCACAGCACAAAATTGGGCACTGGAGCAAAGAGGAGTGGTGTGGAGGCCCTGGCTAAGTATTGACCAATGAGCA

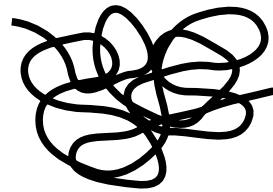


©2010, Illumina Inc. All rights reserved.

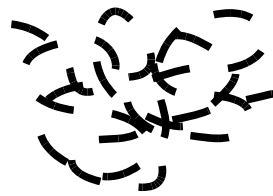
Illumina HiSeq 2000

Variant detection by next-generation sequencing

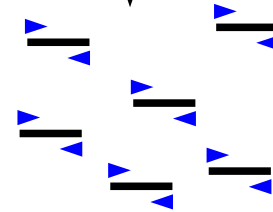
genomic
DNA



fragmented
DNA



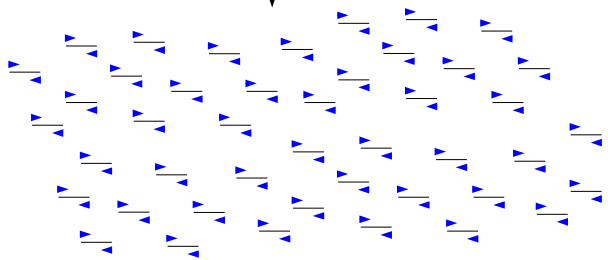
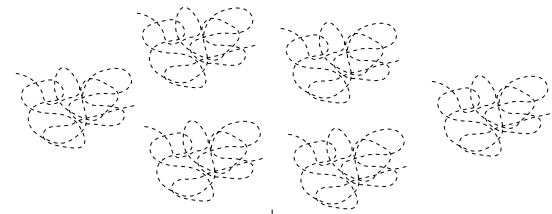
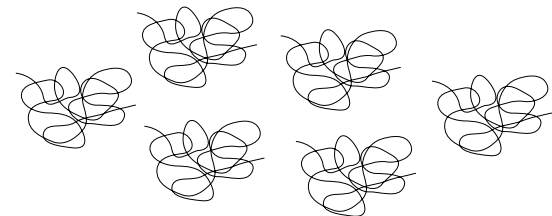
paired-end
sequencing



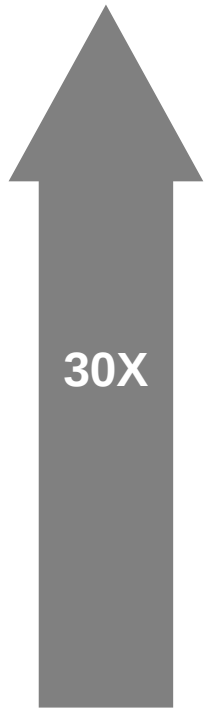
read
mapping

```

GCC ..... GAGAAA GCCACA ..... TGGTCA TGGAGC ..... GCGTGG GCCTGG ..... CAATGA
..... GGTCTG AAAGGC ..... AAAATT CACTGG ..... AGTGGC TGGAGG ..... TTGACC
...TAGATG GTCTGA ..... CA--AC AATTGG ..... AAGAGG TGGCGT ..... AGTATT ACCAAT .....
...ACATGG AGGTCT ..... AGGCCA A--ACAA ..... GCACTG GCAAAG ..... GTGGAG CCAGGC ..... CAATGA
GCCTA ATGGAG ..... GAAAGG ..... ATTGGG ..... AAAGAG AGTGGC ..... GCCTGG AAGTAT ..... TGAGCA
..CTACAT GAGAAT...CACA--A ACTGGA ..... TGGCGT GGCTAA ..... ACGAGC
GCCTACATGGAGGTCTGAGAAAGGCCACAGCACAAAATTGGGCACTGGAGCAAAGAGGAGTGGTGTGGAGGCCCTGGCTAAGTATTGACCAATGAGCA
  
```



Variant detection by next-generation sequencing



```
GCC.....GAGAAA GCCACA.....TGGTCA TGGAGC.....GCGTGG GCCTGG.....CAATGA
.....GGTCTG AAAGGC.....AAAATT CACTGG.....AGTGGC TGGAGG.....TTGACC
...TAGATG GTCTGA.....CA--AC AATTGG.....AAGAGG TGGCGT.....AGTCTT ACCAAT.....
...ACATGG AGGTCT...AGGCCA A--ACAA.....GCACTG GCAAAG.....GTGGAG CCTGGC.....CAATGA
iGCCTA ATGGAG.....G-AAGG          ATTGGG.....AAAGAG AGTGGC.....GCCTGG AAGTAT.....TGAGCA
..CTACAT          GAGAAA...CACA--A          ACTGGA.....TGGCGT          GGCTAA.....ATGAGC
GCC.....GAGAAA GCCACA.....TGGGCA TGGAGC.....GCGTGG GCCTGG.....CAATGA
.....GGTCTG AAAGGC.....AAAATT CACTGG.....AGTGGC TGGAGG.....TTGACC
...TACATG  GCCTGA.....CA--AC AATTGG.....AAGAGG TGGCGT.....AGTATT ACCAAT.....
...ACATGG AGGTCT...AGGCCA A--ACAA.....GCACTG GCAAAG.....GTGGAG C-TGGC.....CAATGA
iGCCTA ATGGAG.....GAAAGG          ATTGGG.....AAAGAG AGTGGC.....GCCTGG AAGTAT.....TGAGCA
..CTACAT          GAGAAA...CACA--A          AGTGGA.....TGGCGT          GGCTAA.....ATGAGC
GCC.....GAGAAA GCCACA.....TGGGCA TGGAGC.....GCGTGG GCCTGG.....CAATGA
.....GGTCTG AAAGGC.....AAAATT CACTGG.....AGTGGC TGGAGG.....TTGACC
...TACATG  GTCTGA.....CA--AC AATTGG.....AAGAGG TGGCGT.....AGTATT ACCAAT.....
...ACATGG AGGTCT...AGGCCA A--ACAA.....GCACTG GCAAAG.....GTGGAG CCTGGC.....CAATGA
iGCCTA ATGGAG.....GAAAGG          ATTGGG.....AAAGAG AGTGGC.....GCCTGG AAGTAT.....TGAGCA
..CTACAT          GAGAAA...CACA--A          ACTGGA.....TGGCGT          GGCTAA.....ATGAGC
GCC.....GAGAAA GCCACA.....TGGGCA TGGAGC.....GCGTGG GCCTGG.....CAATGA
.....GGTCTG AAAGGC.....AAAATT CACTGG.....AGTGGC TGGAGG.....TTGACC
...TACATG  GTCTGA.....CA--AC AATTGG.....AAGAGG TGGCGT.....AGTATT ACCAAT.....
...ACATGG AGGTCT...AGGCCA A--ACAA.....GCACTG GCAAAG.....GTGGAG CCTGGC.....CAATGA
iGCCTA ATGGAG.....GAAAGG          ATTGGG.....AAAGAG AGTGGC.....GCCTGG AAGTAT.....TGAGCA
..CTACAT          GAGAAAT...CACA--A          ACTGGA.....TGGCGT          GGCTAA.....ATGAGC
GCC.....GAGAAA GCCACA.....TGGTCA TGGAGC.....GCGTGG GCCTGG.....CAATGA
.....GGTCTG AAAGGC.....AAAATT CACTGG.....AGTGGC TGGAGG.....TTGACC
...TAGATG  GTCTGA.....CA--AC AATTGG.....AAGAGG TGGCGT.....AGTATT ACCAAT.....
...ACATGG AGGTCT...AGGCCA A--ACAA.....GCACTG GCAAAG.....GTGGAG CCAGGC.....CAATGA
iGCCTA ATGGAG.....GAAAGG          ATTGGG.....AAAGAG AGTGGC.....GCCTGG AAGTAT.....TGAGCA
..CTACAT          GAGAAAT...CACA--A          ACTGGA.....TGGCGT          GGCTAA.....ACGAGC
iGCCTACATGGAGGTCTGAGAAAGGCCACAGCACAAAATTGGGCACTGGAGCAAAGAGGAGTGGTGTGGAGGCCCTGGCTAAGTATTGACCAATGAGCA
```

read
mapping

Next-Generation Sequencing (NGS)

Throughput

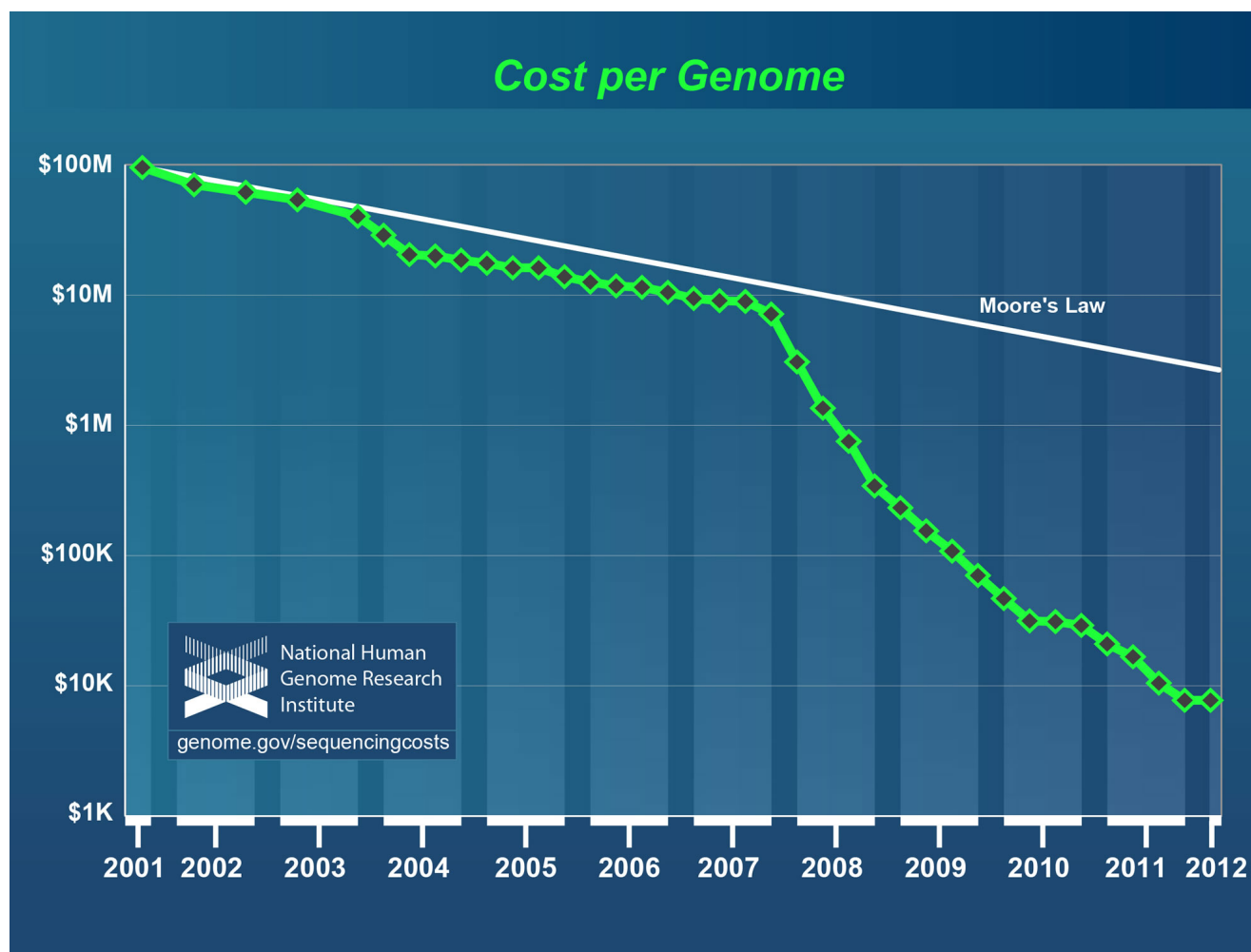


Illumina HiSeq 2000 run

- 10 days runtime
- ~6 billion 100bp paired-end reads
- ~30x coverage of 6 human genome
- ~220Gb of (compressed) raw data

Next-Generation Sequencing (NGS)

Reduction in DNA sequencing costs



Next-Generation Sequencing (NGS)

Informatics Challenges

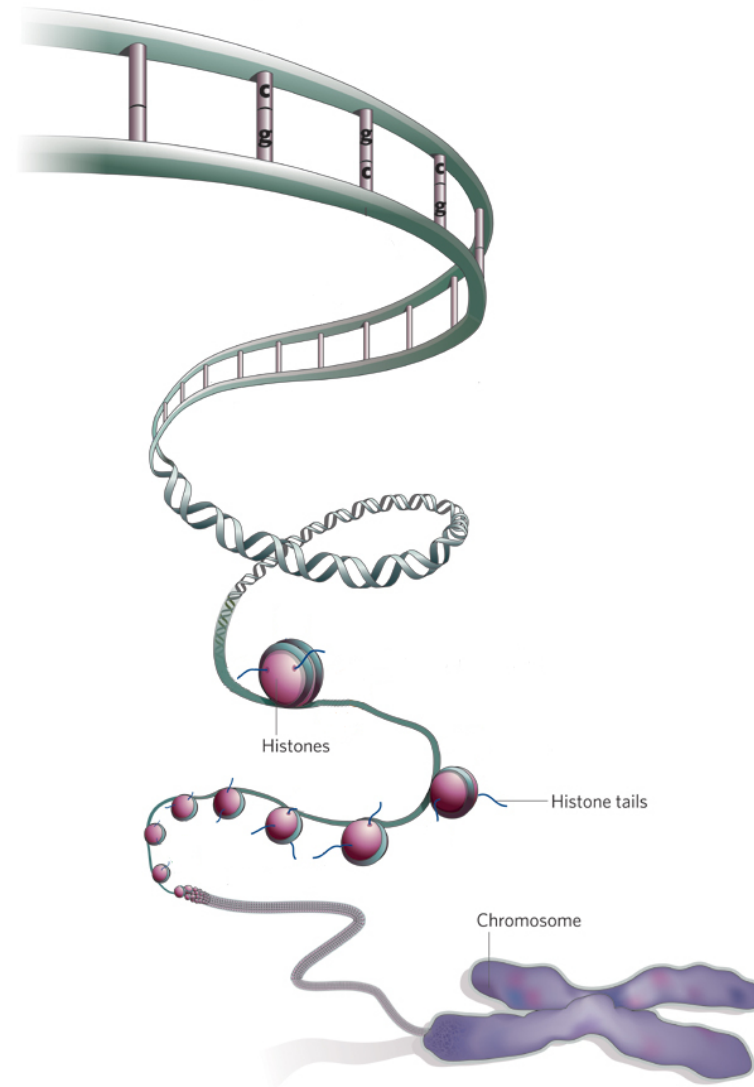
- **data storage & I/O**
 - ~7TB raw data per sequencer per year
- **mapping algorithm efficiency**
 - 1st generation: hash-based
 - 2nd generation: Burrows-Wheeler transform
 - 10 times faster at similar sensitivity
- **CPU time**
 - 720 CPU hours (30 days) for ~1 billion reads (30x coverage of the human genome)

Next-Generation Sequencing (NGS)

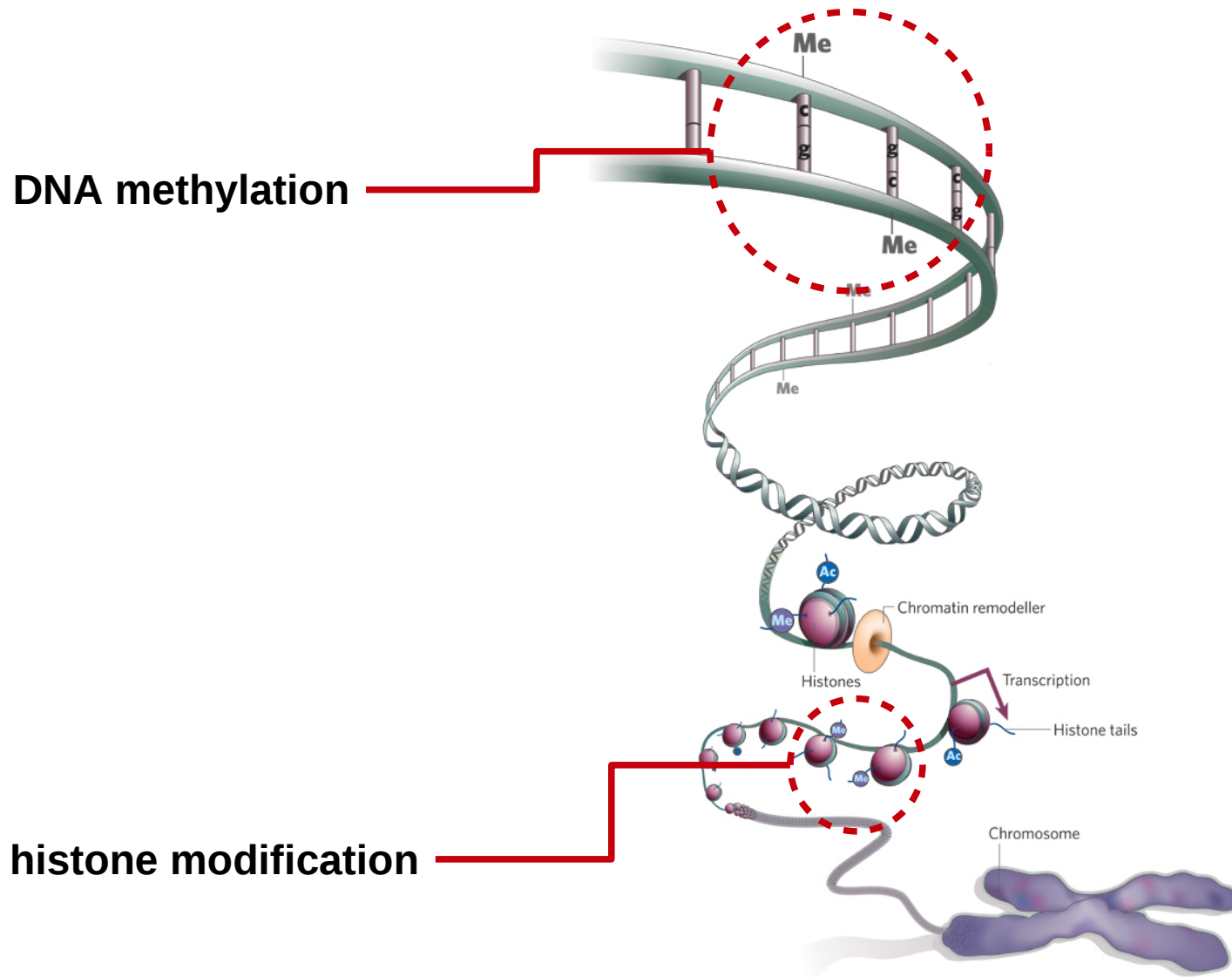
NGS data analysis using Imperial College HPC service

- dedicated access to 112 CPU cores on PC cluster
- access to SGI Altix UV system with 256 CPU cores and up to 3TB of shared memory
- 60TB storage

From genome to epigenome



From genome to epigenome

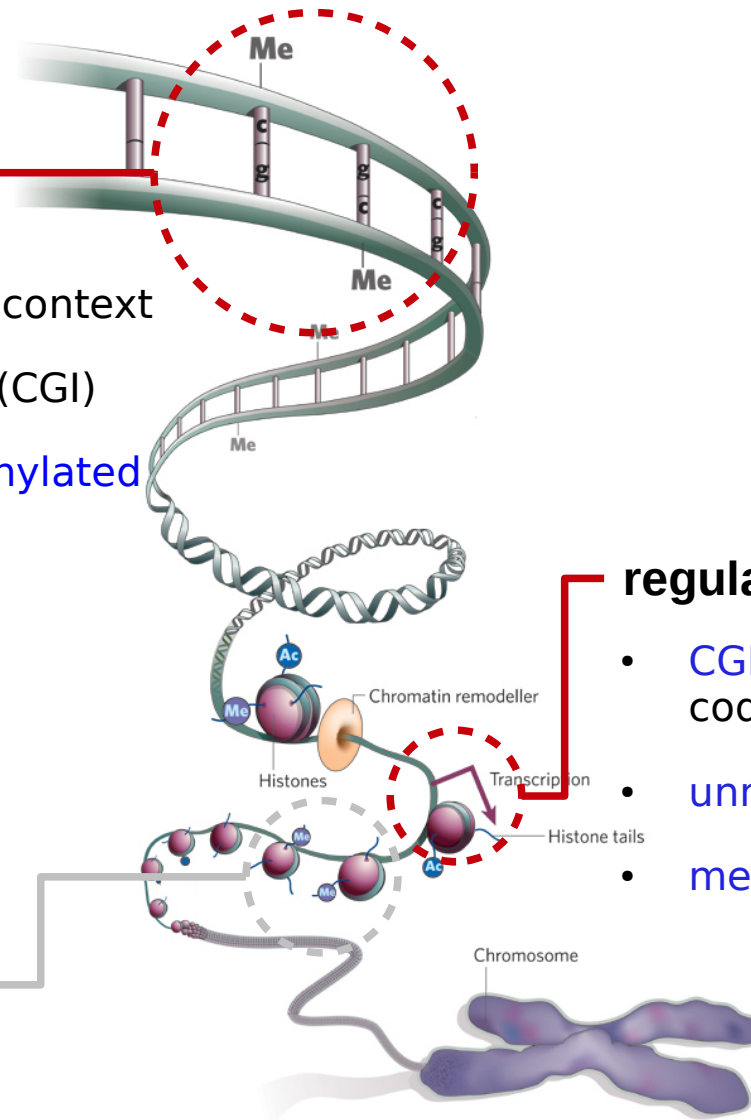


DNA methylation & gene expression

DNA methylation

- occurs in CpG dinucleotides context
- CpGs cluster in CpG Islands (CGI)
- genome globally highly methylated
- CGIs mostly unmethylated

histone modification



regulation of transcription

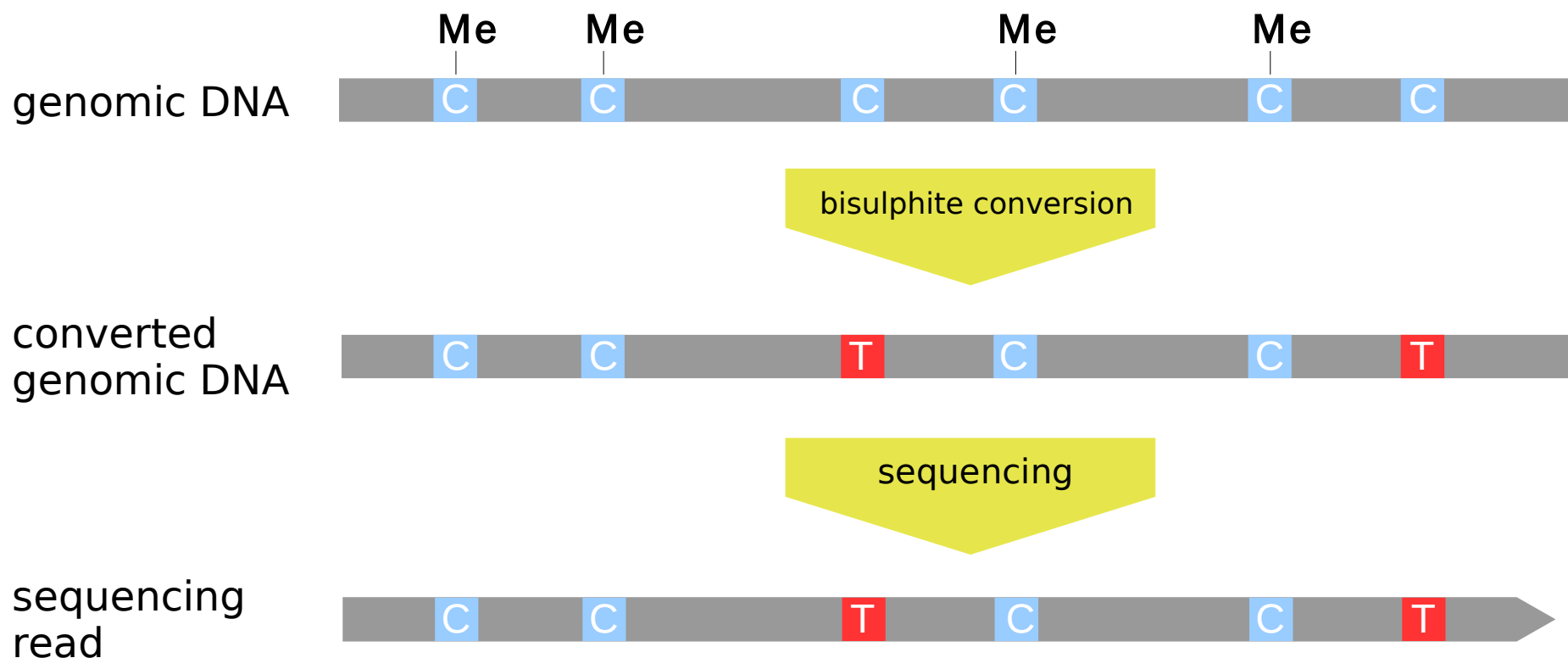
- CGIs found in 50% of protein coding gene promoters
- unmethylated promoters active
- methylated promoters repressed

Genetic control of inter-individual variation in DNA methylation

Aims

- measure DNA methylation levels at single nucleotide resolution in two rat inbred strains (Brown Norway & Spontaneously Hypertensive Rat)
- identify inter-strain differences in DNA methylation
- determine inheritance and genetic regulators of DNA methylation

Measuring DNA methylation by bisulphite sequencing



Measuring DNA methylation by bisulphite sequencing



bisulphite
sequencing
reads



Measuring DNA methylation by bisulphite sequencing



bisulphite
sequencing
reads



Measuring DNA methylation by bisulphite sequencing



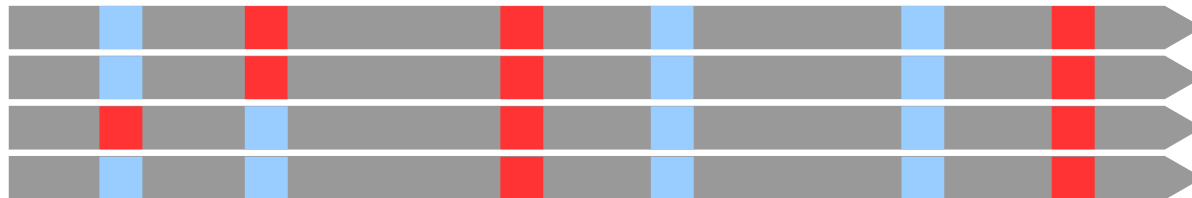
bisulphite
sequencing
reads



Measuring DNA methylation by bisulphite sequencing



bisulphite
sequencing
reads



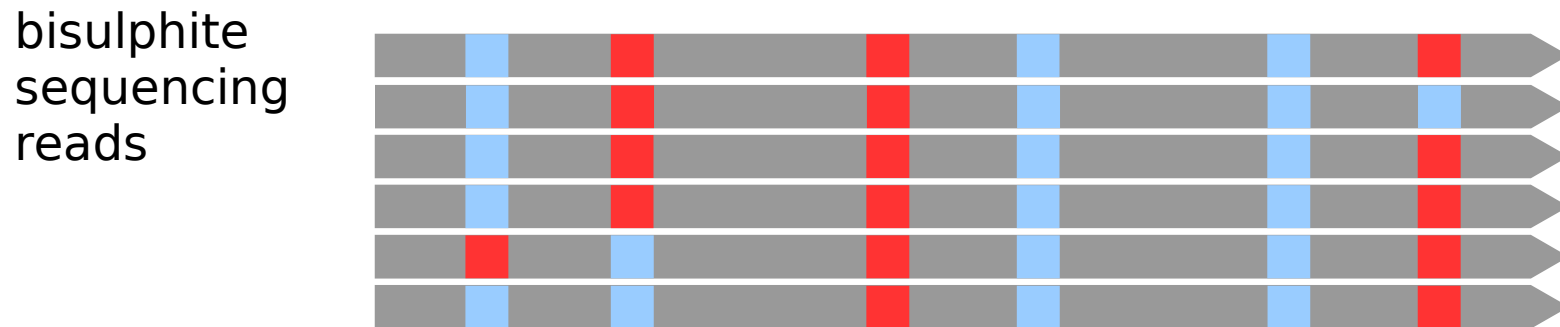
Measuring DNA methylation by bisulphite sequencing



bisulphite
sequencing
reads



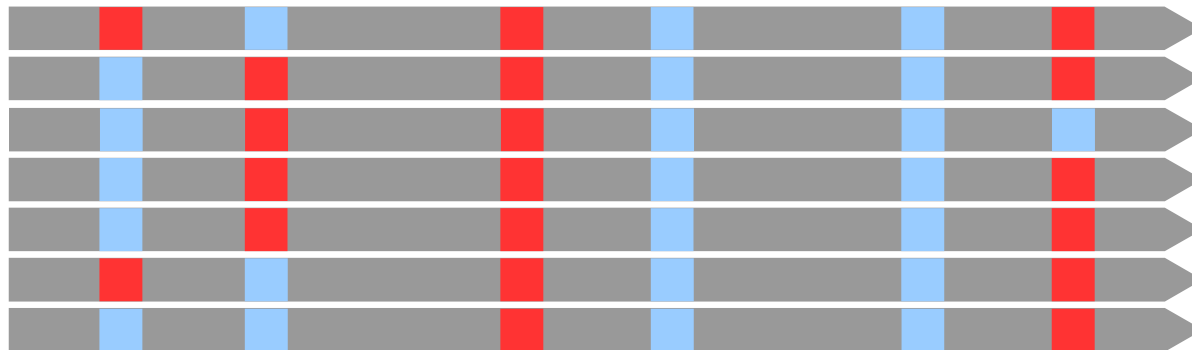
Measuring DNA methylation by bisulphite sequencing



Measuring DNA methylation by bisulphite sequencing



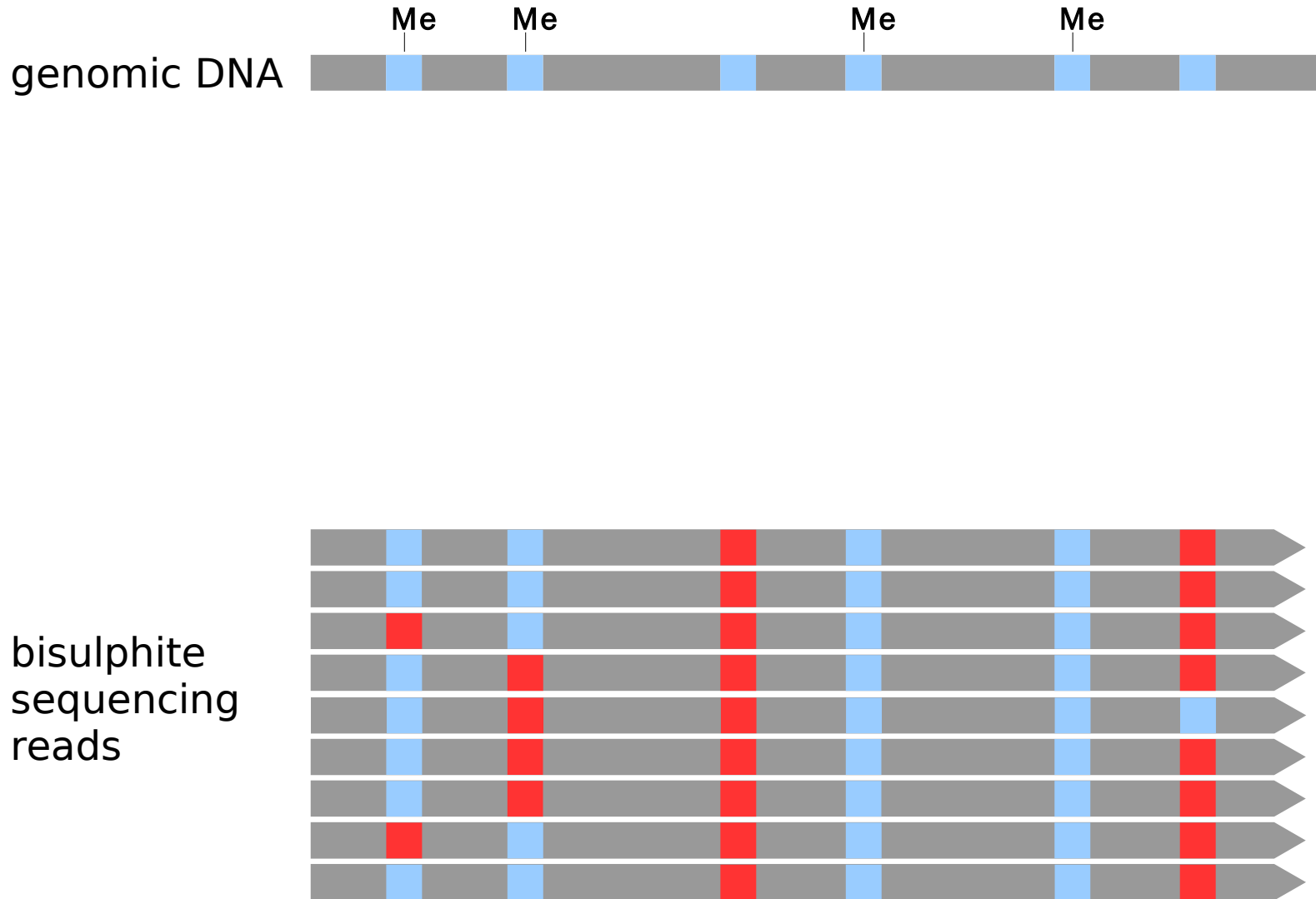
bisulphite
sequencing
reads



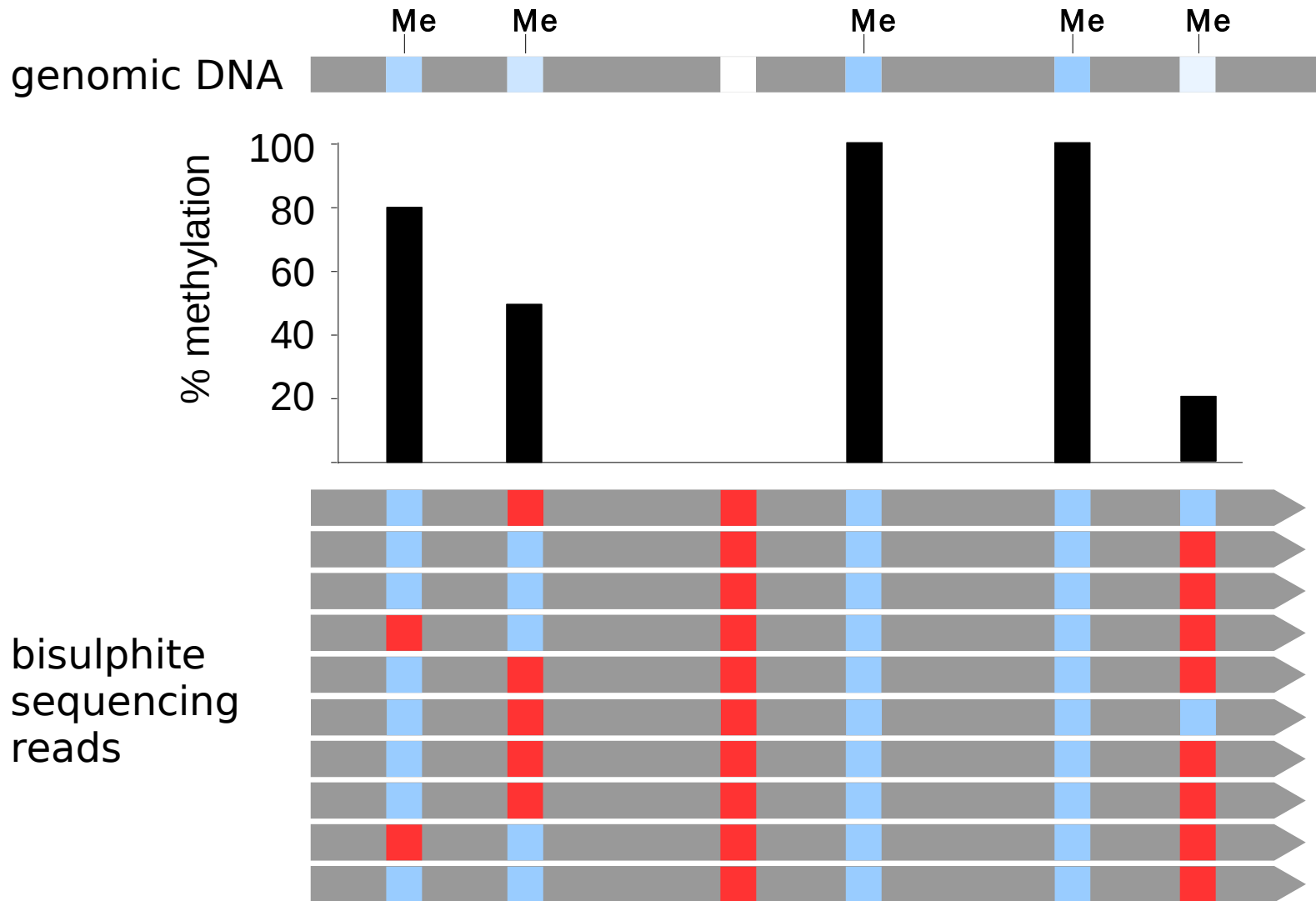
Measuring DNA methylation by bisulphite sequencing



Measuring DNA methylation by bisulphite sequencing



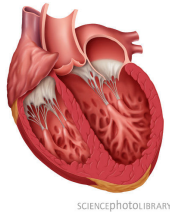
Measuring DNA methylation by bisulphite sequencing



Sequencing the SHR & BN methylomes



- Illumina paired-end sequencing
- 100 bp read length
- 4 lanes per sample



- left ventricle of the heart



- SHR/Ola, BN-Lx, reciprocal F1 crosses
- 4 biological replicates
- six weeks old

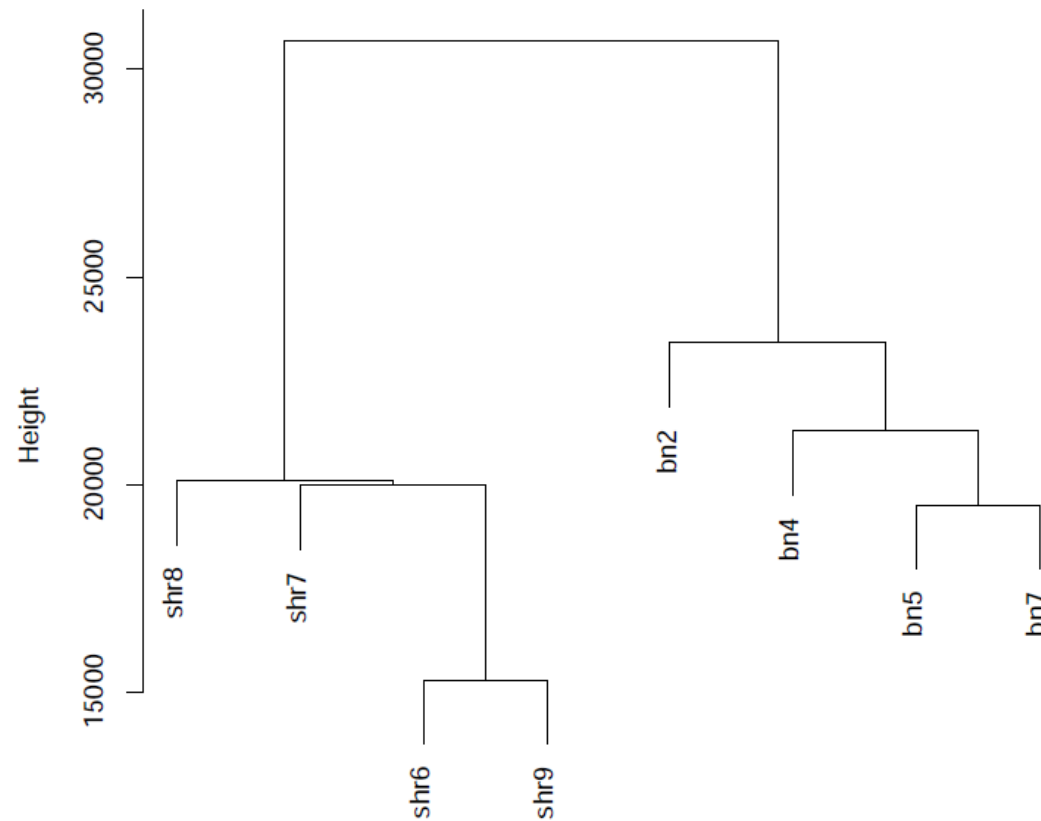
Sequencing the rat methylome

sample	reads	mapped reads	depth of coverage	depth of coverage after filtering*
	[million]	[million]	[%]	[x-fold]
BN2	258	230	89	7
BN4	413	361	87	13
BN5	454	404	89	11
BN7	370	340	92	12
BN total	1,495	1,336	89	42
SHR6	497	448	91	15
SHR7	452	398	89	12
SHR8	421	369	90	11
SHR9	391	356	91	12
SHR total	1,762	1,571	89	50

* filtering by read clonality, uniqueness, mapping quality and strand specificity

Inter- vs intra-strain differences in CpG methylation

Variability in CpG methylation greater between than within inbred strains

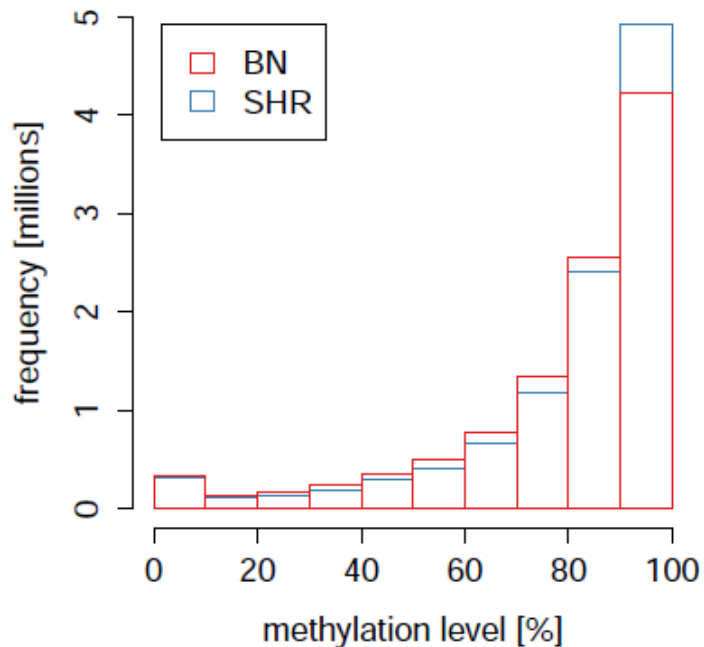


Distance measure: Euclidean
Clustering method: Ward

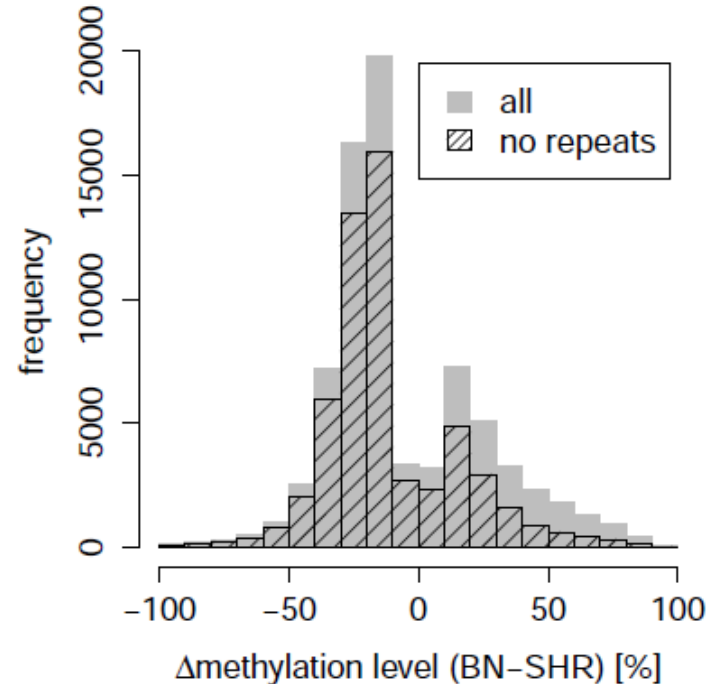
Distribution of methylation levels and inter-strain differences in methylation

analysis of 10.6 million CpG dinucleotides identified
77,088 differentially methylated CpGs (<1%)

The rat genome is
globally highly methylated

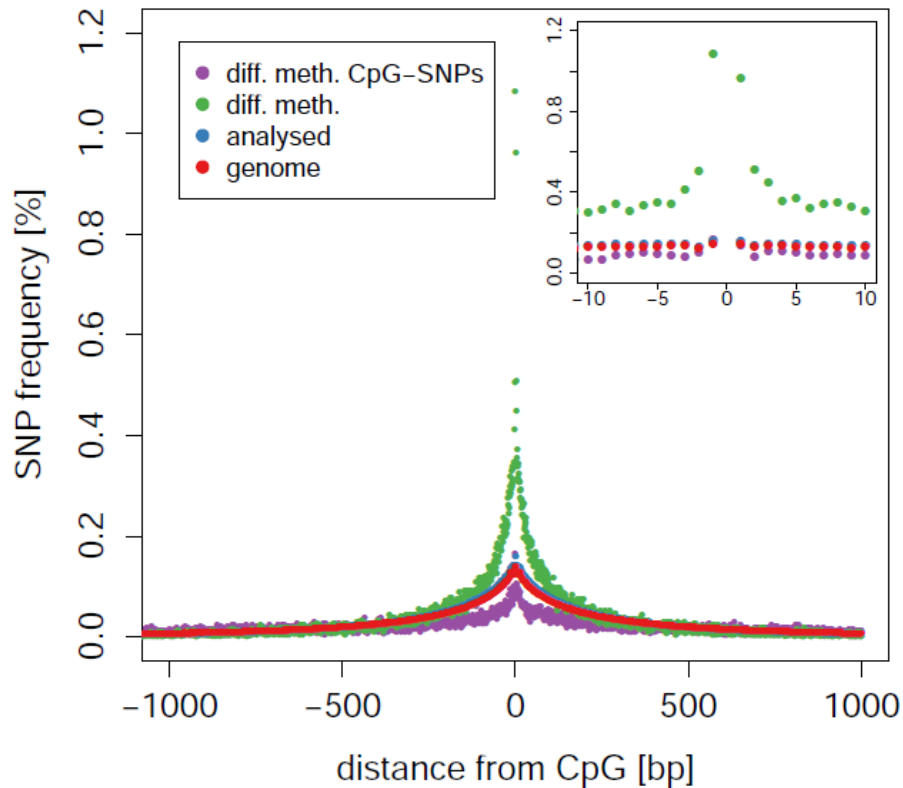


The majority of methylation
differences are small

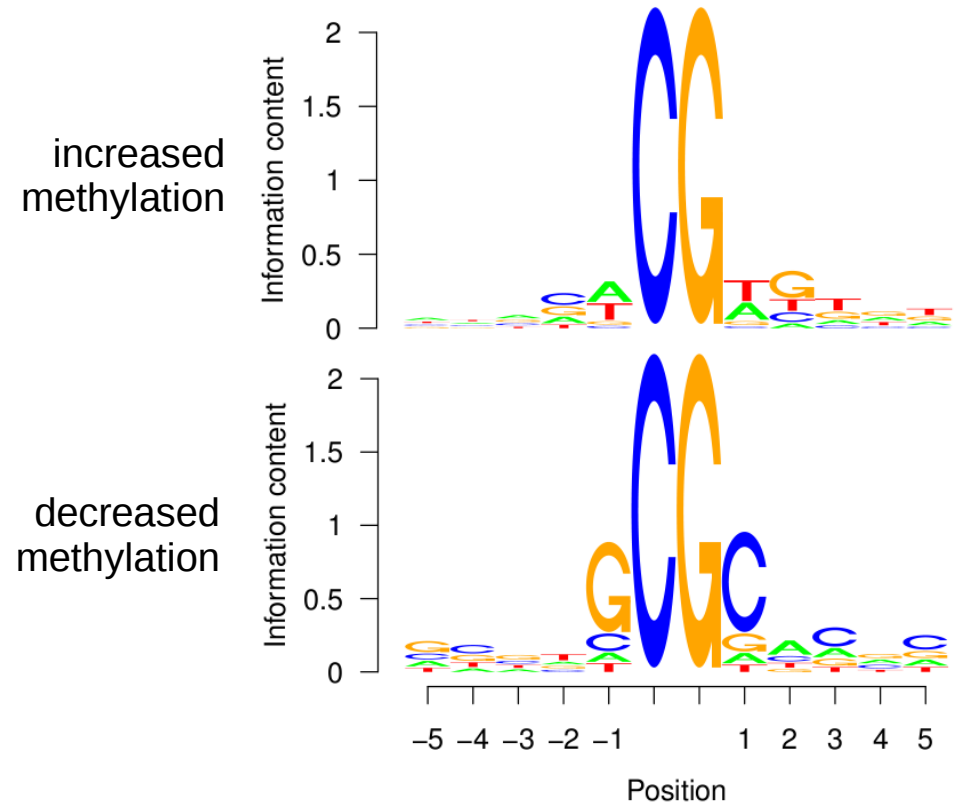


Sequence variation at differentially methylated CpGs

Sequence variation is higher in the proximity of differentially methylated CpGs

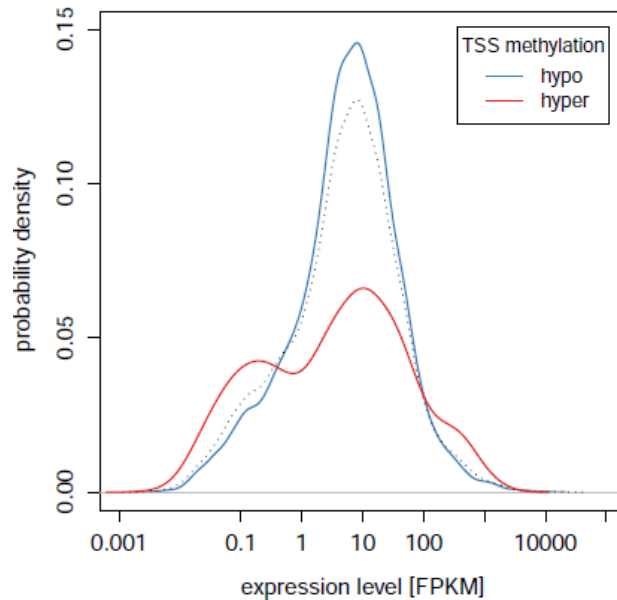


Increased and decreased methylation is associated with specific sequence motifs

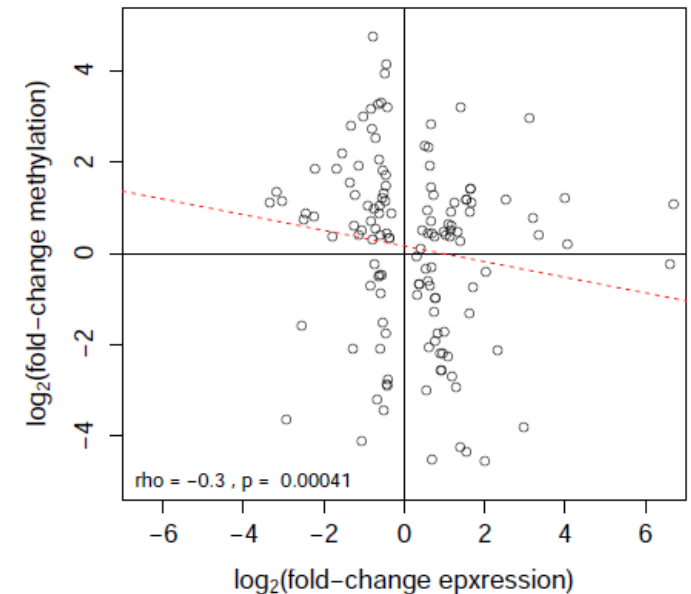


Methylation and gene expression

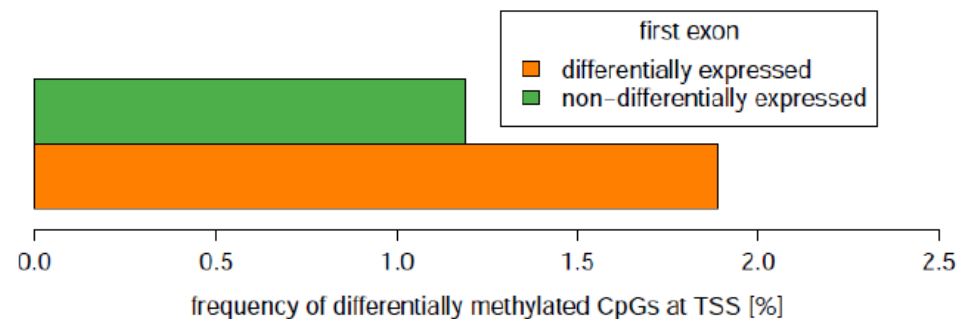
Genes with highly and lowly methylated promoters have different distribution of expression levels



Methylation and gene expression differences are negatively correlated



Promoters of differentially expressed genes are enriched in differentially methylated CpGs



Genetic control of inter-individual variation in DNA methylation

Summary

- generated whole-genome methylation profiles at single nucleotide resolution in two rat strains (SHR & BN)
- quantified methylation differences between SHR & BN
 - inter-strain methylation differences are greater than intra-strain differences
 - majority of methylation differences are small
 - SNP frequency in proximity of differential methylation CpGs is increase
 - increased/decreased methylation associated with specific sequence motifs
- integrated methylation and gene expression profiles
 - clear but complex association of differential methylation and differential expression

Acknowledgments

Physiological Genomics & Medicine Group

Michelle Johnson
(Post-doctoral Research Associate)

Klio Maratou
(Post-doctoral Research Associate)

Prashant Shrivastava
(Post-doctoral Research Associate)

Prof. Timothy J. Aitman
(Group Head)

Cardiac Regeneration Group National Heart & Lung Institute

Prof. Nadia A. Rosenthal
(Head of the Heart Science
Centre)

MRC CSC Genomics Laboratory

Laurence Game
(Facility Head)

Adam Giess
(Bioinformatician)

Korinne Northwood
(Research Assistant)

Imperial College High Performance Computing

Simon Burbidge
(HPC Co-ordination Manager)

