

# Cancer Epigenetics Study Using Next-Generation Sequencing Data

**July 29, 2010**  
**Big Data For Science**

**Sun Kim and Heejooon Chae**  
*School of Informatics and Computing*  
*Center for Bioinformatics Research*  
*Indiana University, Bloomington, Indiana, USA*

# Overview of The Talk

- Background on epigenomics and DNA methylation
- OSU-IU Center for Cancer Systems Biology
- Mapping sequence reads
- Data
- BioVLAB-mCpG

# Part I: Epigenomics and DNA Methylation

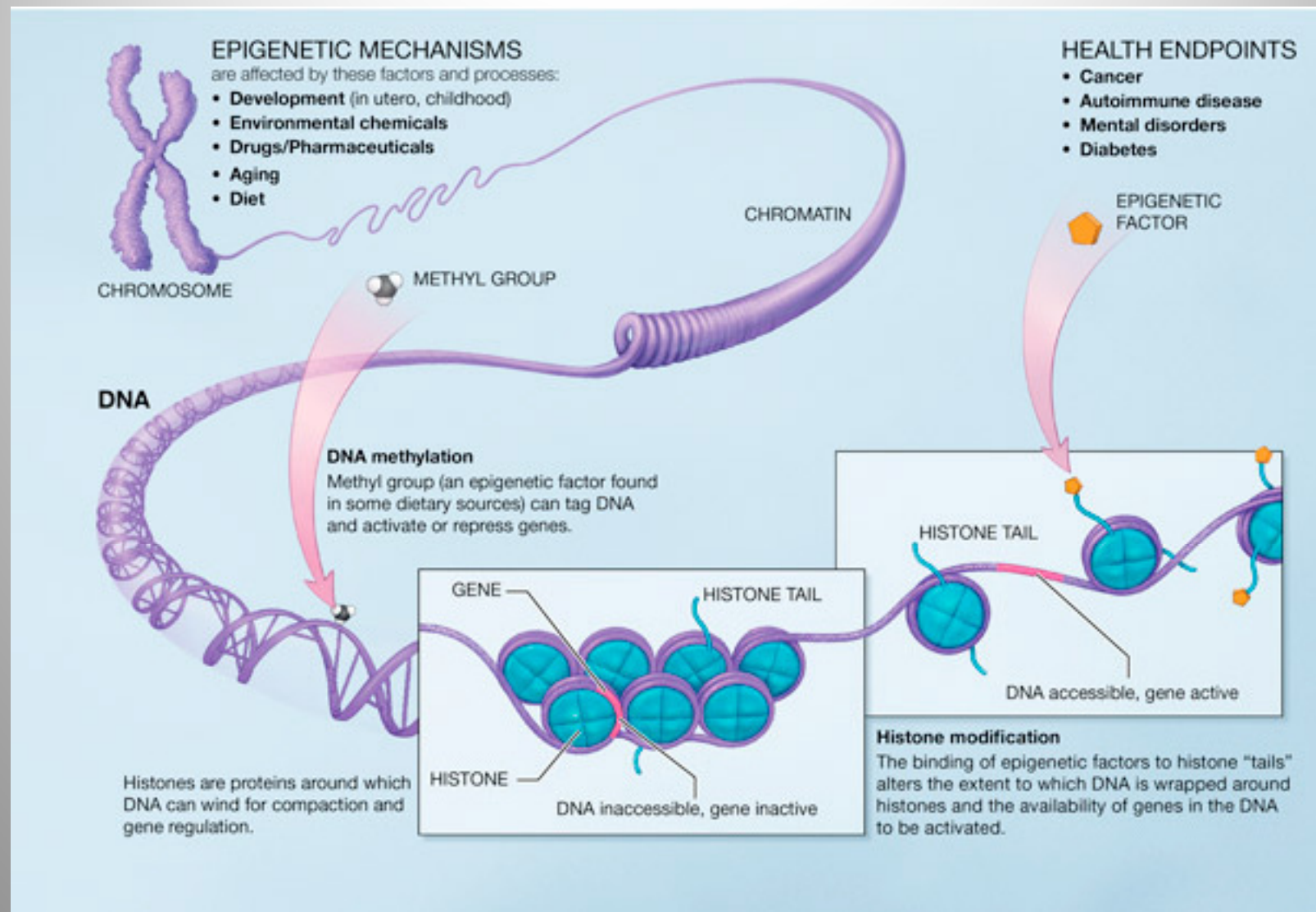
-- Sun Kim group at IU --

3

# Epigenetics

- Epigenetics is the study of **heritable** changes in gene function that occur without a change in DNA sequence.
- Summarizes mechanisms and phenomena that affect the phenotype of a cell or an organism without affecting the genotype.
- Modifications of DNA (cytosine methylation) and proteins (histones) define the epigenetic profile.
- *Epigenomics* is the study of these epigenetic changes on a genome-wide scale.

*This slide is from Ken Nephew at IU.*

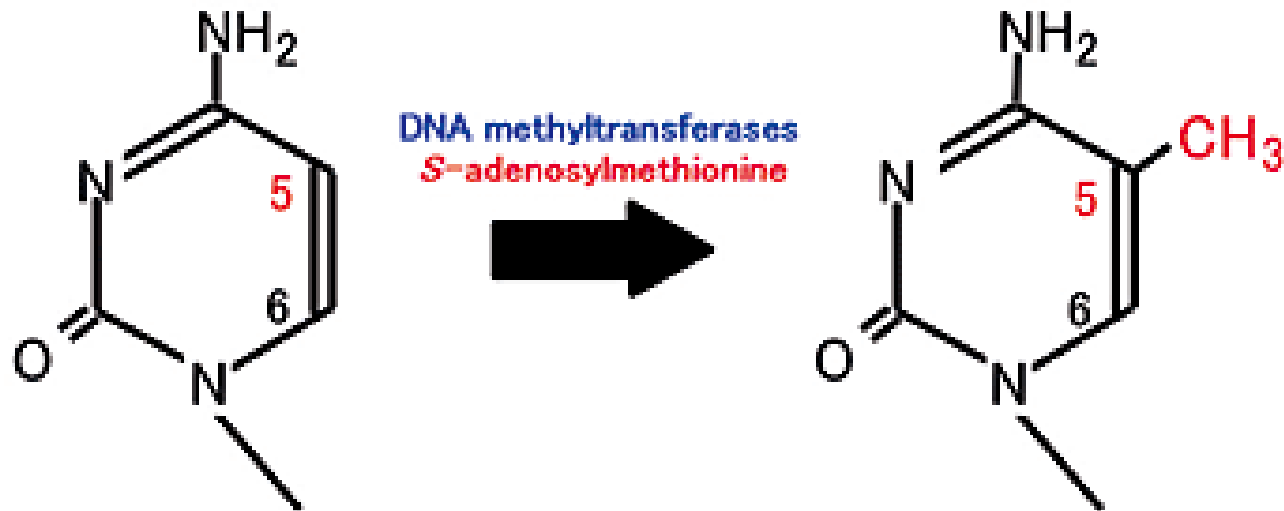


<http://nihroadmap.nih.gov/epigenomics/epigeneticmechanisms.a>

-- Sun Kim group at IU --

# DNA Methylation

## Methylation of Cytosine



5'—CpG—3'

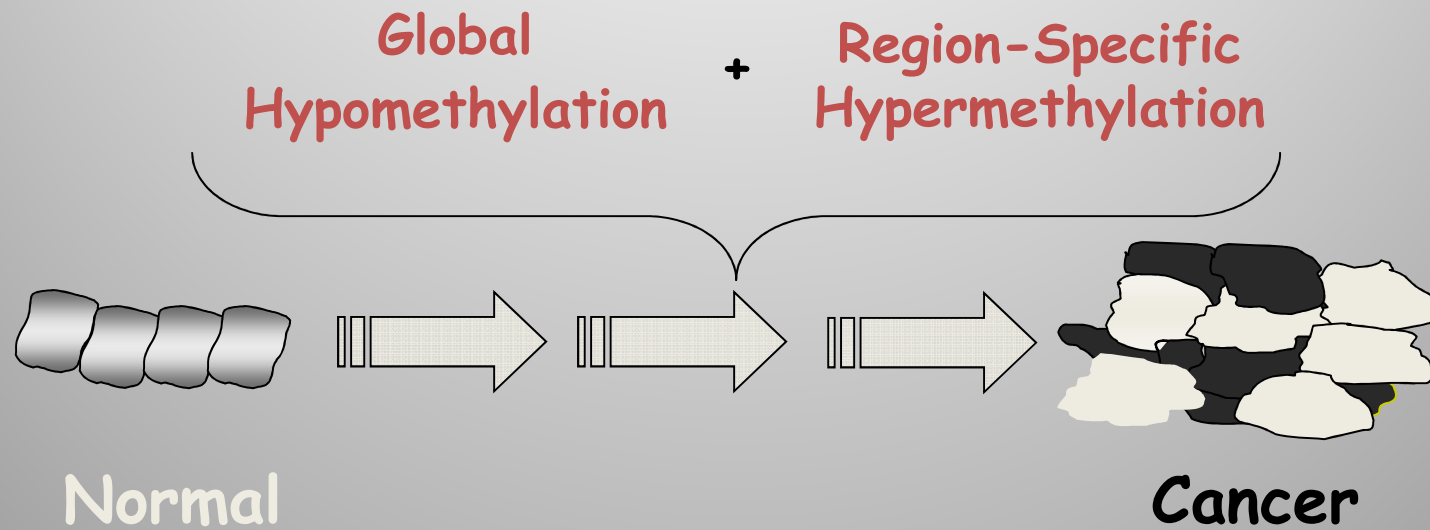
3'—GpC—5'

# **Normal Cellular Functions Regulated by Epigenetic Mechanisms**

- **Correct organization of chromatin**
  - Controls active and inactive states of embryonic and somatic cell-  
Epigenetic components contribute to plasticity and stability during development.
  - Involved in maintenance of differentiated cells.
- **Specific DNA methylation patterns, chromatin modifications**
  - Controls gene- and tissue-specific epigenetic patterns.
- **Genomic imprinting**- Essential for development
- **Silencing of repetitive elements**
  - Maintains chromatin order, proper gene expression patterns
- **X chromosome inactivation**- Balances gene expression

*This slide is from Ken Neptew at IU.*

## Progressive Accumulation of DNA Methylation in Cancer



Accumulation of Genetic and Epigenetic Abnormalities

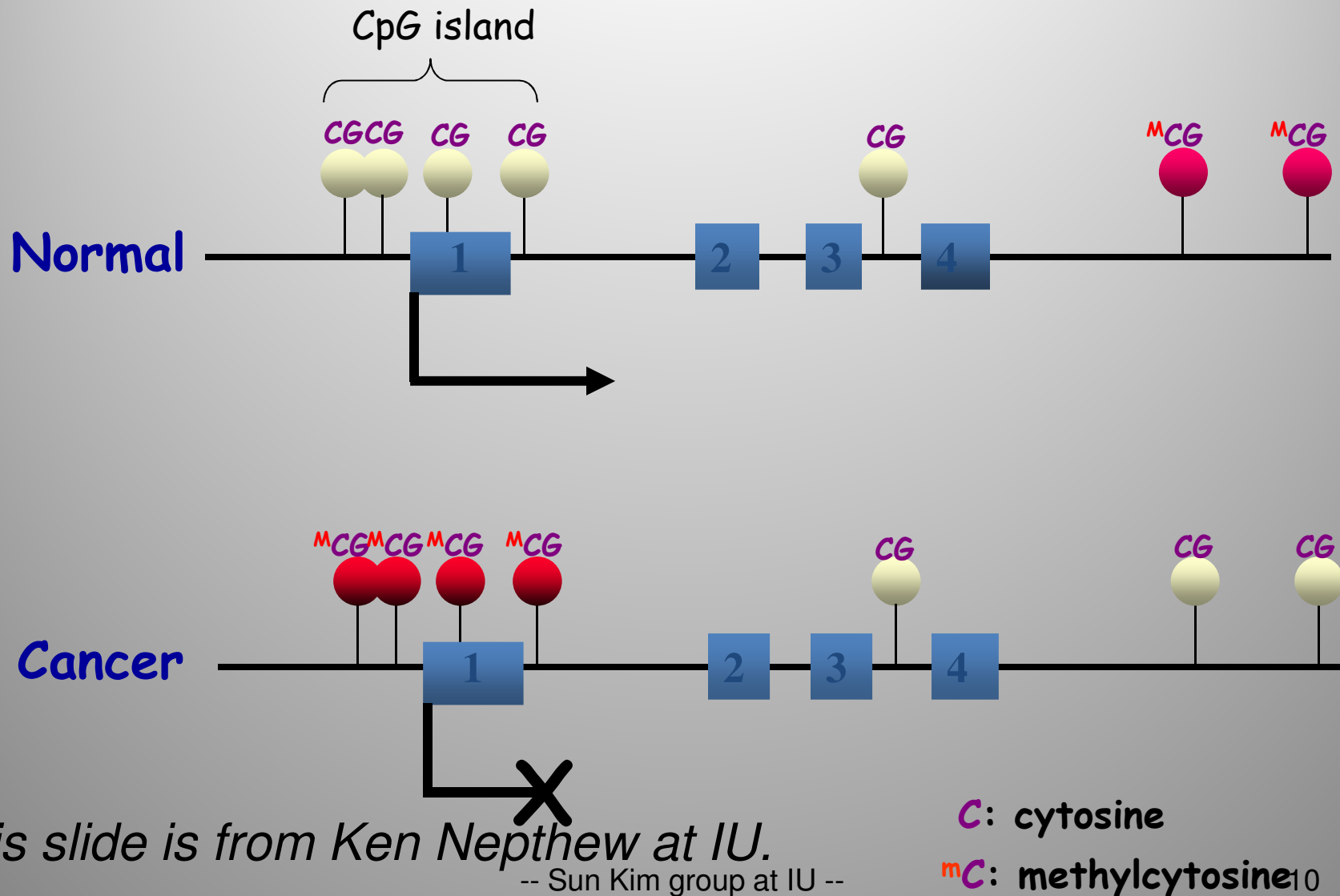
*This slide is from Ken Nephew at IU.*



# CpG Islands

- CpG island: a cluster of CpG residues often found near gene promoters (sequences ~1000 base pairs in length with a GC content of over 60%)
- ~29,000 CpG islands in human genome (~60% of all genes are associated with CpG islands)
- Most CpG islands are unmethyated in normal cells.

# DNA Methylation and Gene Silencing in Cancer Cells



*This slide is from Ken Nephew at IU.*

-- Sun Kim group at IU --

# Histone modifications: Histone Code

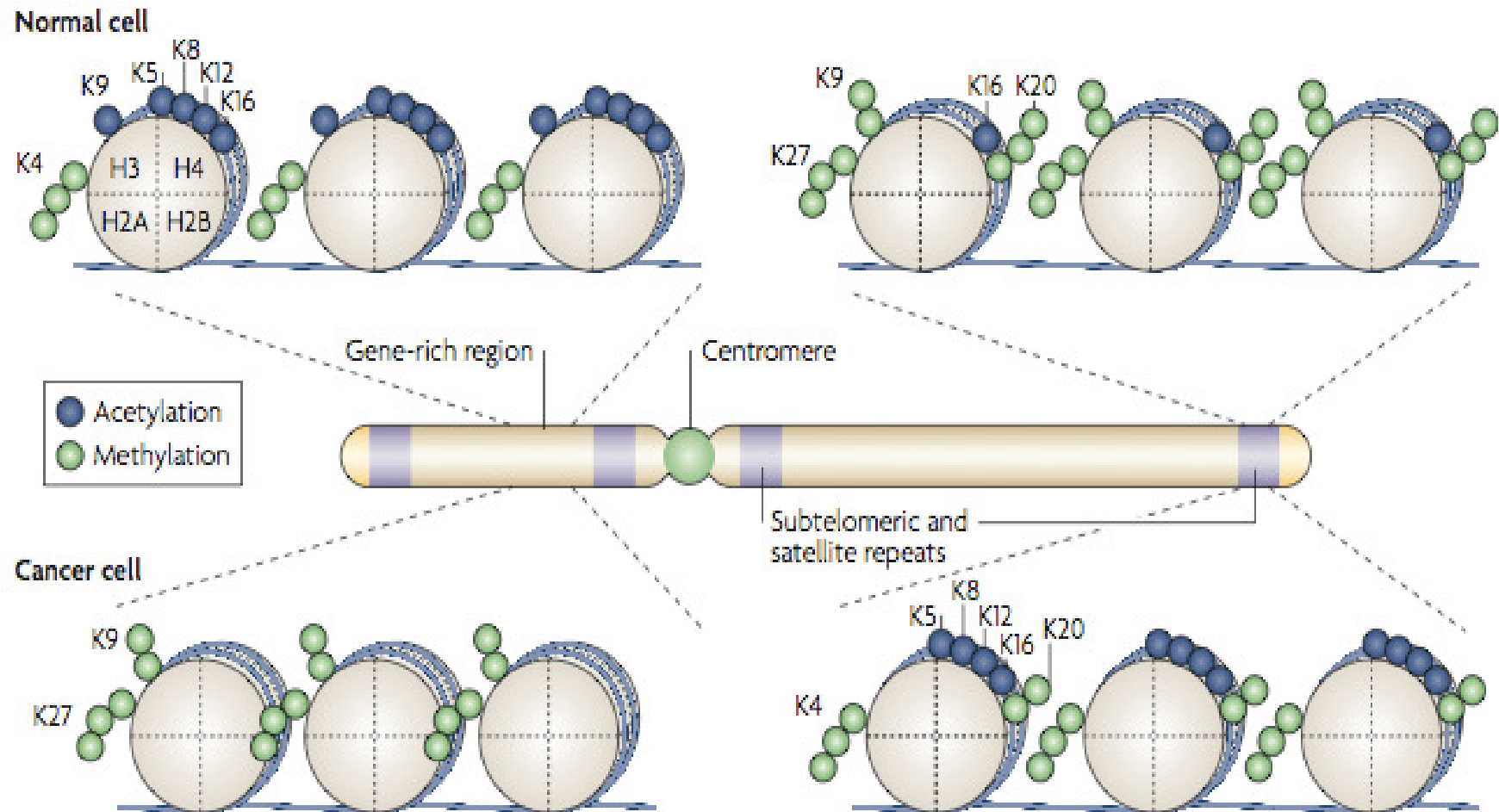
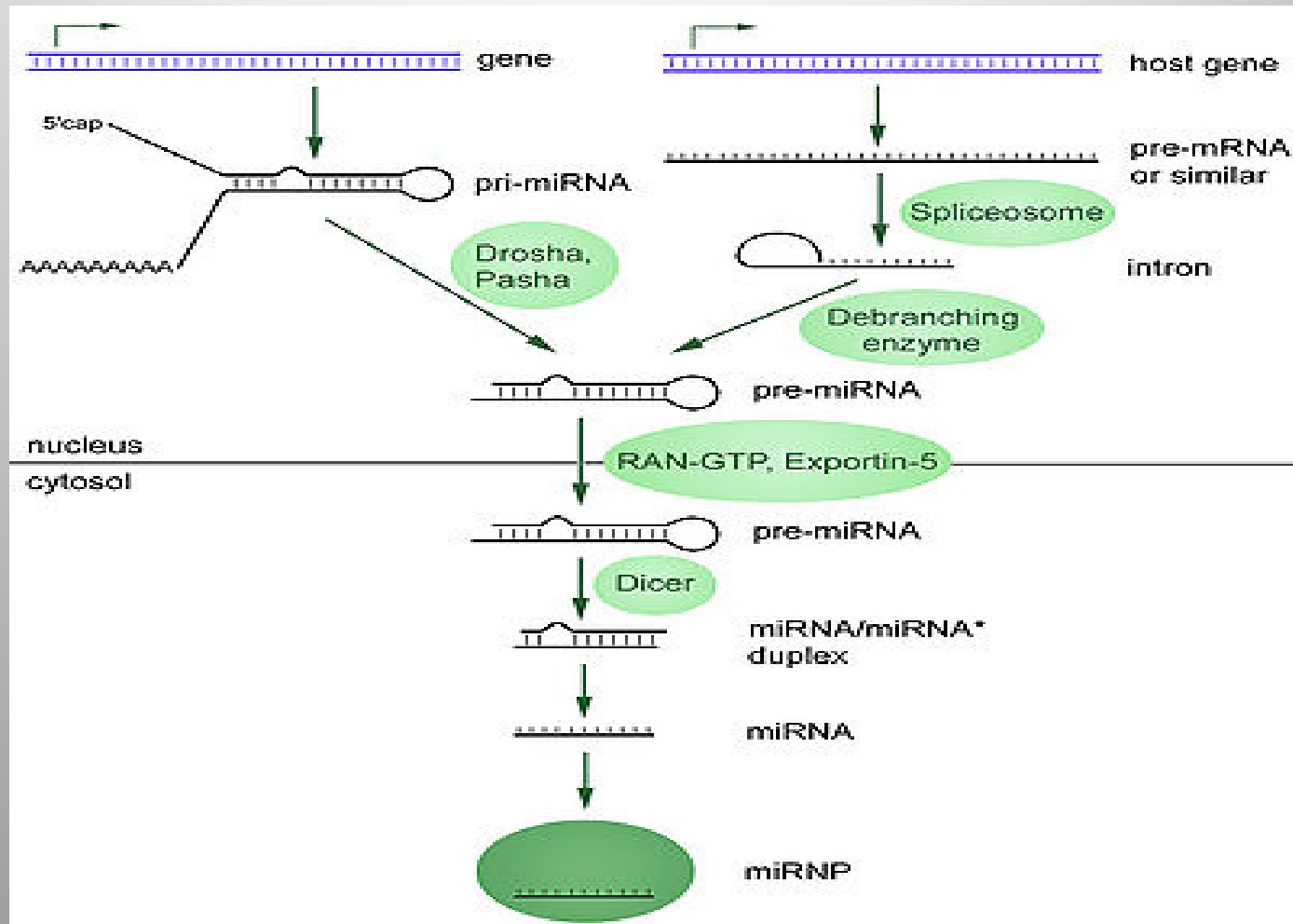


Figure 4 | **Histone-modification maps for a typical chromosome in normal and cancer cells.** Nucleosomal arrays

Nature Reviews Genetics 8, 286-298 (April 2007)

# MicroRNA



<http://en.wikipedia.org/wiki/MicroRNA>

## **PART 2: OSU-IU Center for Cancer Systems Biology**

-- Sun Kim group at IU --

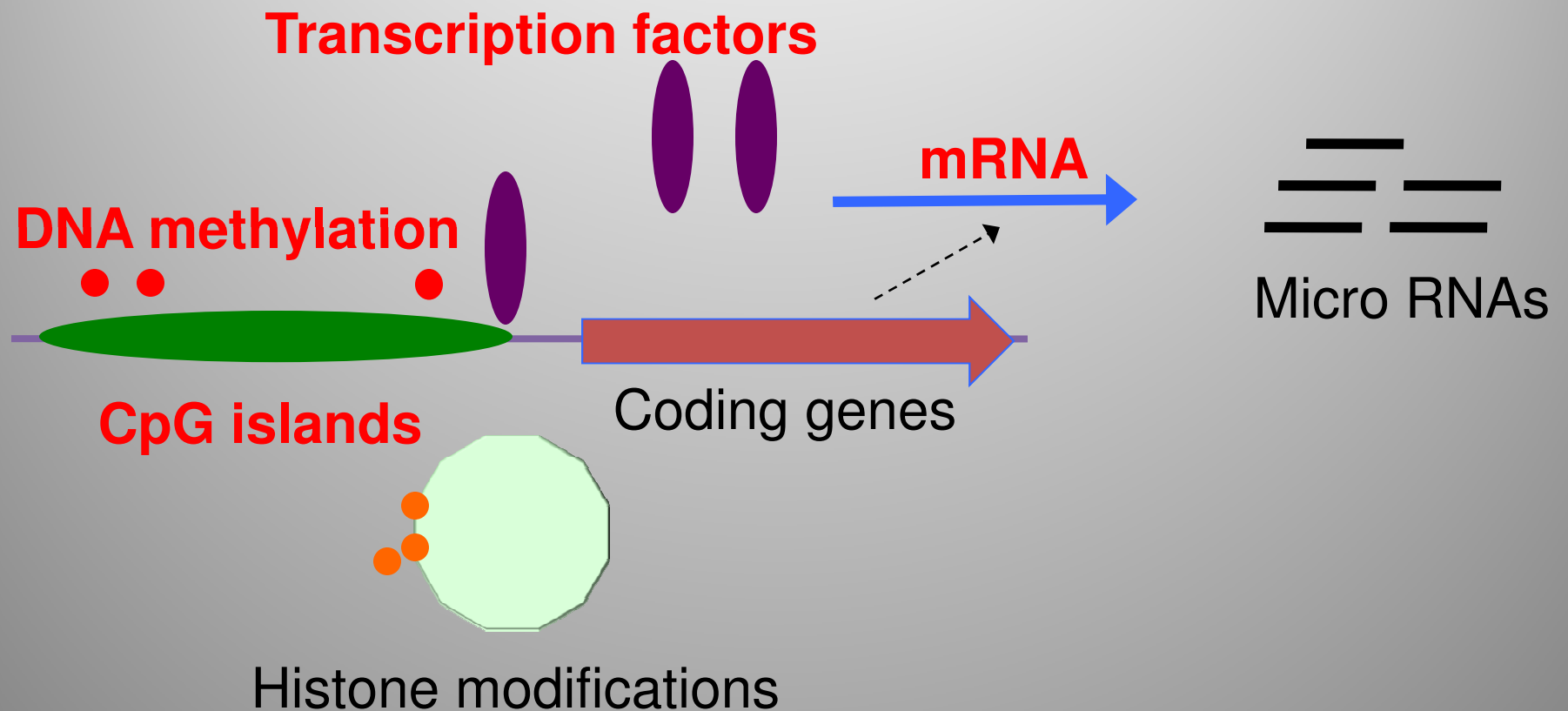
# **OSU-IU Integrated Cancer Biology Program (ICBP) Center**

- The Integrative Cancer Biology Program (<http://icbp.nci.nih.gov/>) is a program launched by US National Cancer Institute in 2004.
- OSU-IU ICBP Center aims to characterize the role of epigenomics in the development of drug resistance in human cancer for a period of 2004 – 2015.

# Drug Resistance in Human Cancer

- The OSU-IU Center has been investigating the mechanism of developing drug resistance in breast, prostate, and ovarian cancer.
- In particular, we are interested in investigating changes in *epigenetic mechanisms* in terms of gene regulation and pathway activation while in transition to a hormone-/chemo-sensitive to ***a hormone-/chemo-insensitive phenotype*** in cancer.

# DNA Methylation vs. Transcription Factor





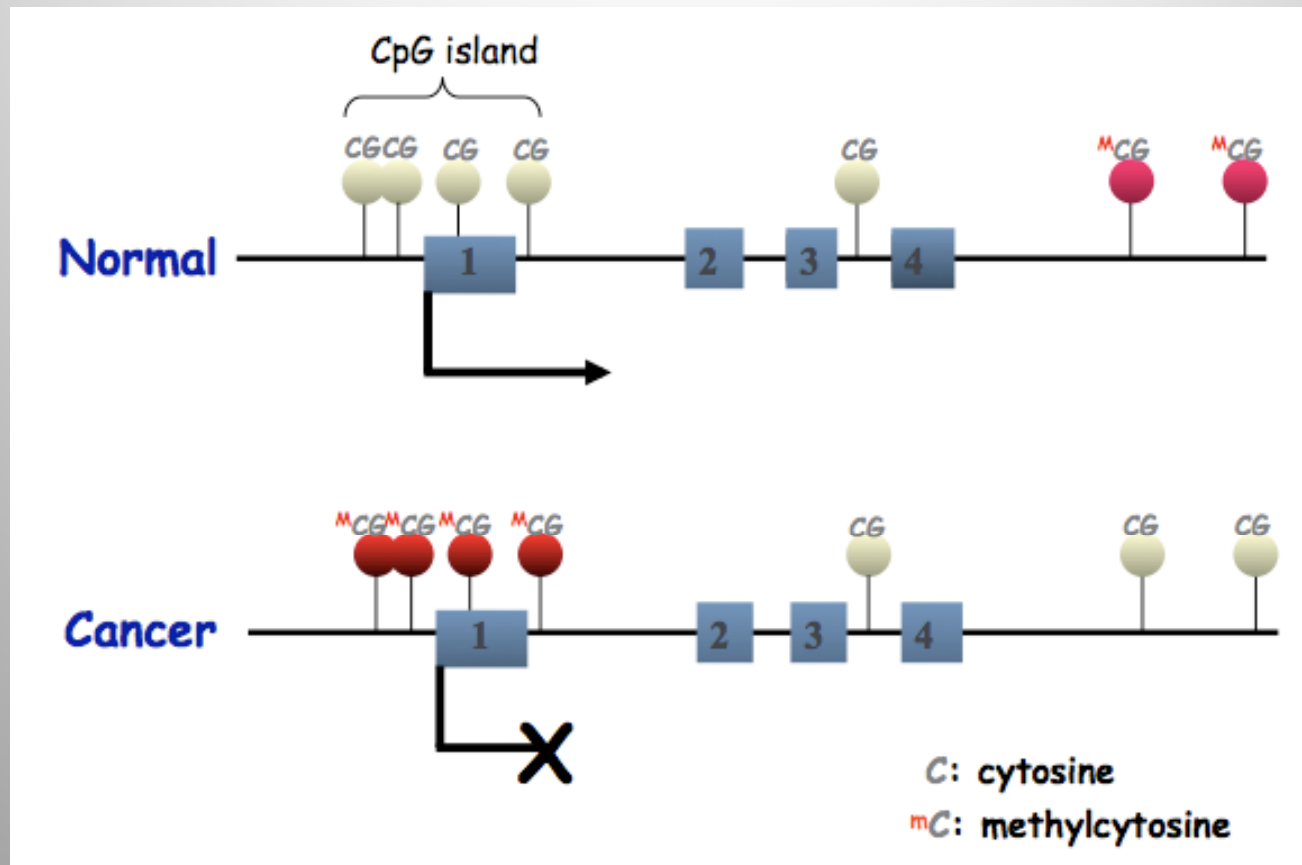
# 6 Methylome Projects

- To investigate the effect of DNA methylation in drug-resistance cancer phenotype, we sequence and study 6 cell lines:
  1. Breast cancer: 2 cell lines before and after drug resistance phenotype.
  2. Prostate cancer: 2 cell lines before and after drug resistance phenotype.
  3. Ovarian cancer: 2 cell lines before and after drug resistance phenotype.

# Basic Data Analysis

- Comparing methylation difference in two cell lines (e.g., before and after drug-resistance phenotype).
- Integrated analysis with histone modification, microRNA, gene expression, and phenotypes.

# Comparative Analysis of Methylation in Two Cell Lines



- Promotor methylation analysis and expression of downstream genes.
- Promotor methylation and transcription factors and their binding sites.
- Intergenic methylation and alternative splicing.
- Methylation in non-CpG context.

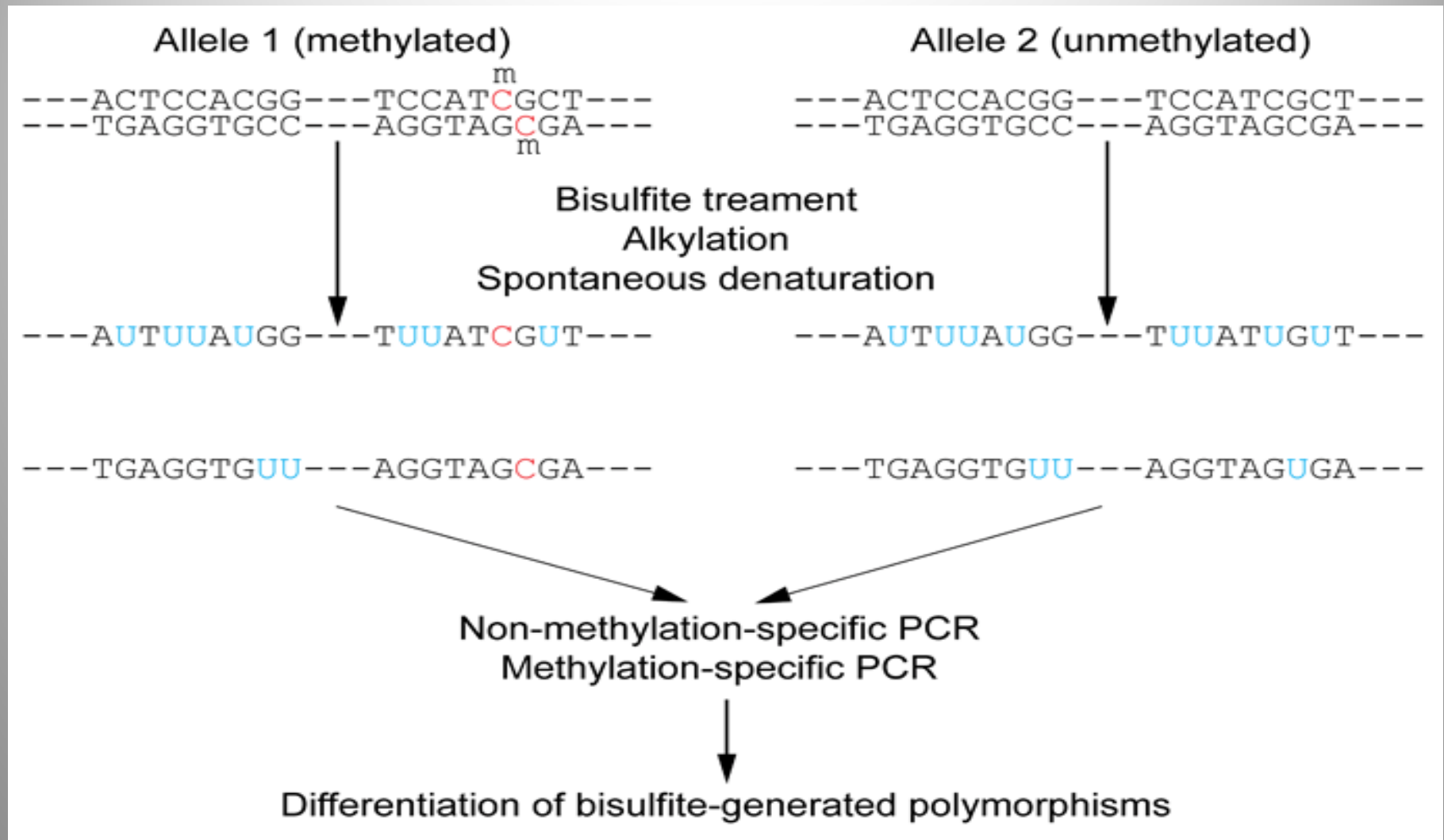
-- Sun Kim group at IU --

# PART 3: Sequence read mapping

-- Sun Kim group at IU --

20

# Bisulfite Sequencing to Identify Methylated Cytosines



[http://en.wikipedia.org/wiki/Bisulfite\\_sequencing](http://en.wikipedia.org/wiki/Bisulfite_sequencing)

# Challenges in Mapping Sequence Reads from Bisulfite Treated DNA

- A lot of reads should be mapped: several hundred millions to several billions.
- To know which cytosines are methylated, we need to sequence bisulfite treated DNA. This results in dealing with sequences of alphabet size 3, thus it takes more time.

# Example of Bisulfite Sequencing

ADAM12_CLL_F_283	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAATGGTGCGC
ADAM12_FL_F_141	TTTGCATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAATGGTGCGC
ADAM12_CLL_F_182	TTTGGATAGTTTGTTTATTTATTGTAATGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_CLL_F_271	TTTGGATAGTTTGTTTATTTATTGTAATGGTTAAGGTTGGTTTGTGTTGGAACGGTGTGC
ADAM12_MCL_F_203	TTTGGATAGTTTGTTTATTTATTGTAATGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGCT
ADAM12_MCL_F_322	TTTGGATAGTTTGTTTATTTATTGTAATGGTTAAGGTTGGTTTGTGTCAGAACGGTGTGT
ADAM12_Normal_F_445	TTTGCATAGTTTGTTTATTTATTGTAATGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_ALL_F_233	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_MCL_F_154	TTTGGATAGTTTGTTTATTTATTGCAACCGGTTAAGGTTGGTTTGTGTTATAACGGTGTGC
ADAM12_CLL_F_161	TTTGGACAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_CLL_F_301	TTCCGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAATGGTGTGC
ADAM12_MCL_F_327	TTTGG--AAGTTGTTTATTTATTT--ACGGTTAAGGTTGGTTTGTGTTAGAATGGTGTGC
ADAM12_ALL_F_270	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_FL_F_139	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGG--CGC
ADAM12_CLL_F_375	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_ALL_F_254	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAATGGTGTGC
ADAM12_ALL_F_477	TTCCGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAATGGTGTGC
ADAM12_ALL_F_431	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_FL_F_36	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_CLL_F_373	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_CLL_F_355	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGTGTGC
ADAM12_CLL_F_145	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGCGCGC
ADAM12_CLL_F_367	TTTGGATAGTTTGTTTATTTATTGTAACCGGTTAAGGTTGGTTTGTGTTAGAACGGCGCGC

Methylation status of ADAM12 gene promotor region:  
courtesy by Huidong Shi at Medical College of Georgia.

# Performance Comparison of Mapping Algorithms

Program	Benchmark data		Human genome	
	Time	# mapped	Time	# mapped
GNUMap	47.9 s	<b>71 262</b>	985 m 14 s	<b>7 739 321</b>
Bowtie	7.0 s	62 298	14 m 43 s	6 699 526
SOAP	11.7 s	62 208	32 m 20 s	6 764 050
MAQ	46.5 s	62 208	*3488 m 28 s	6 764 054
Slider	16 m 31 s*	58 551	Crashed	Crashed
SeqMap	81.2 s	56 326	1703 m 04 s	5 455 538
Novocraft	24.4 s	56 238	*920 m 25 s	5 306 782
RMAP	9.2 s	1202	*295 m 54 s	3 447 086

Bold values show that GNUMAP achieves the best performance.

From *Bioinformatics*. 2010 Jan 1;26(1):38-45

-- Sun Kim group at IU --



# PART 4: Data

-- Sun Kim group at IU --

# Two data sets

- 6 methylome data sets from our center
- 2 cell line data from

Lister R, Pelizzola M, Downen RH, Hawkins RD, Hon G, Tonti-Filippini J, Nery JR, Lee L, Ye Z, Ngo QM, Edsall L, Antosiewicz-Bourget J, Stewart R, Ruotti V, Millar AH, Thomson JA, Ren B, Ecker JR. Human DNA methylomes at base resolution show widespread epigenomic differences. *Nature*. 2009 Nov 19;462(7271):315-2

# Data and Runtime Estimation

Estimate of CPU & storage requirement

Aligner	Data set	# of read/ cell line	CPU hrs/ cell line	# of cell line	Total CPU hrs	Temporal Max Disk
GNUMAP	OSU	80 million	950	6	5700	1000GB
	Nature	4 billion	47500	2	95000	13TB
	TBD	4 billion	47500	2	95000	13TB
Bowtie	OSU	80 million	150	6	900	1000GB
	Nature	4 billion	750	2	1500	13TB
	TBD	4 billion	750	2	1500	13TB
Sub total					199000	13TB
<u>Bisulfite treatment reads need 4times run per each read</u>					X4	
<b>Total</b>					<b>796000</b>	<b>13TB</b>

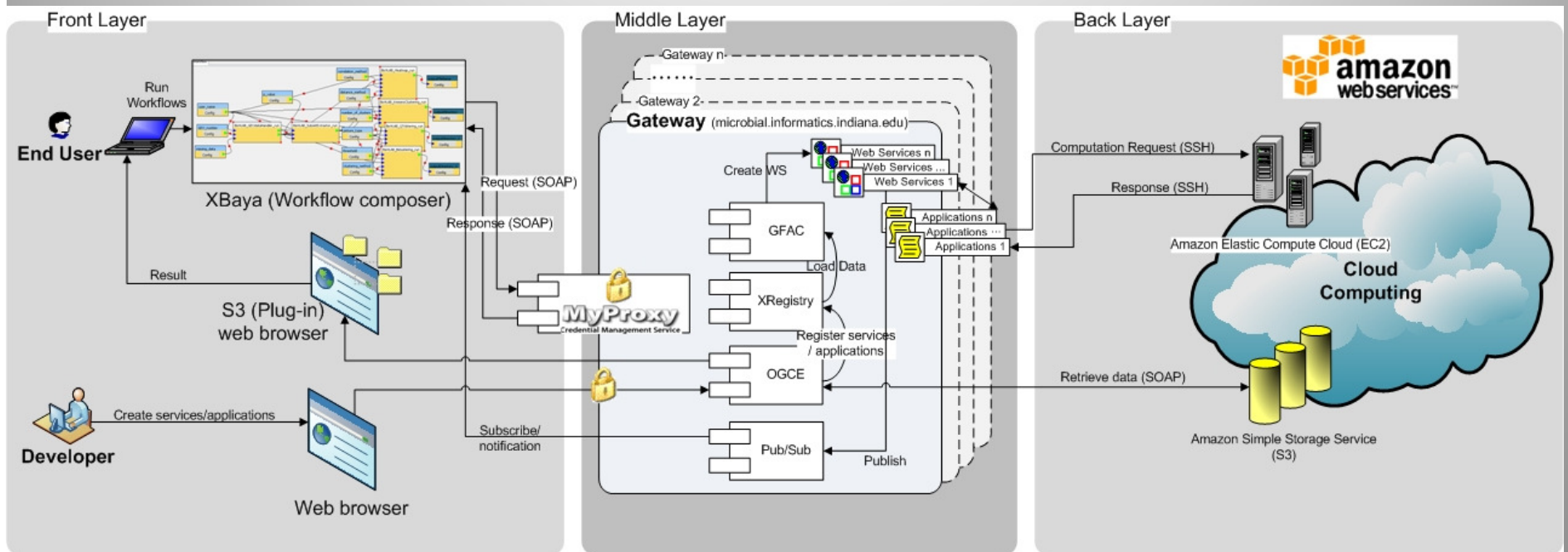
# PART 5: BioVLAB-mCpG

-- Sun Kim group at IU --

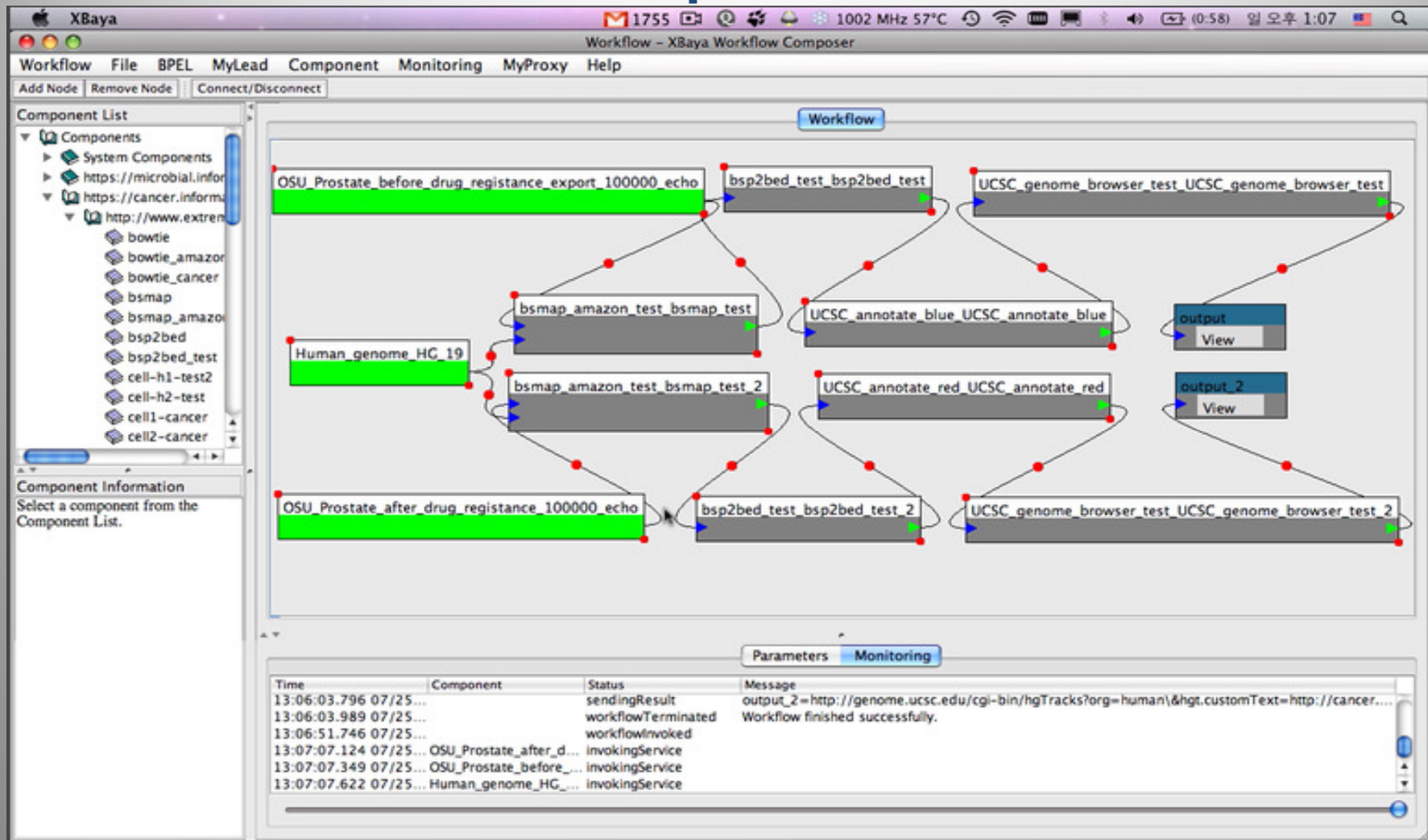
# BioVLAB: Motivation

- We have developed a computational infrastructure, called BioVLAB, for the analysis of molecular biology data utilizing Amazon Cloud Computing (or any high performance computing machines) and a graphical workflow composer, XBay.
- Easy to perform computational analysis:
  1. Set up an account
  2. Download a precomposed workflow
  3. (Modify workflow if needed: application-specific cloud)
  4. Run it

# BioVLAB Architecture



# BioVLAB-mCpG Screenshots

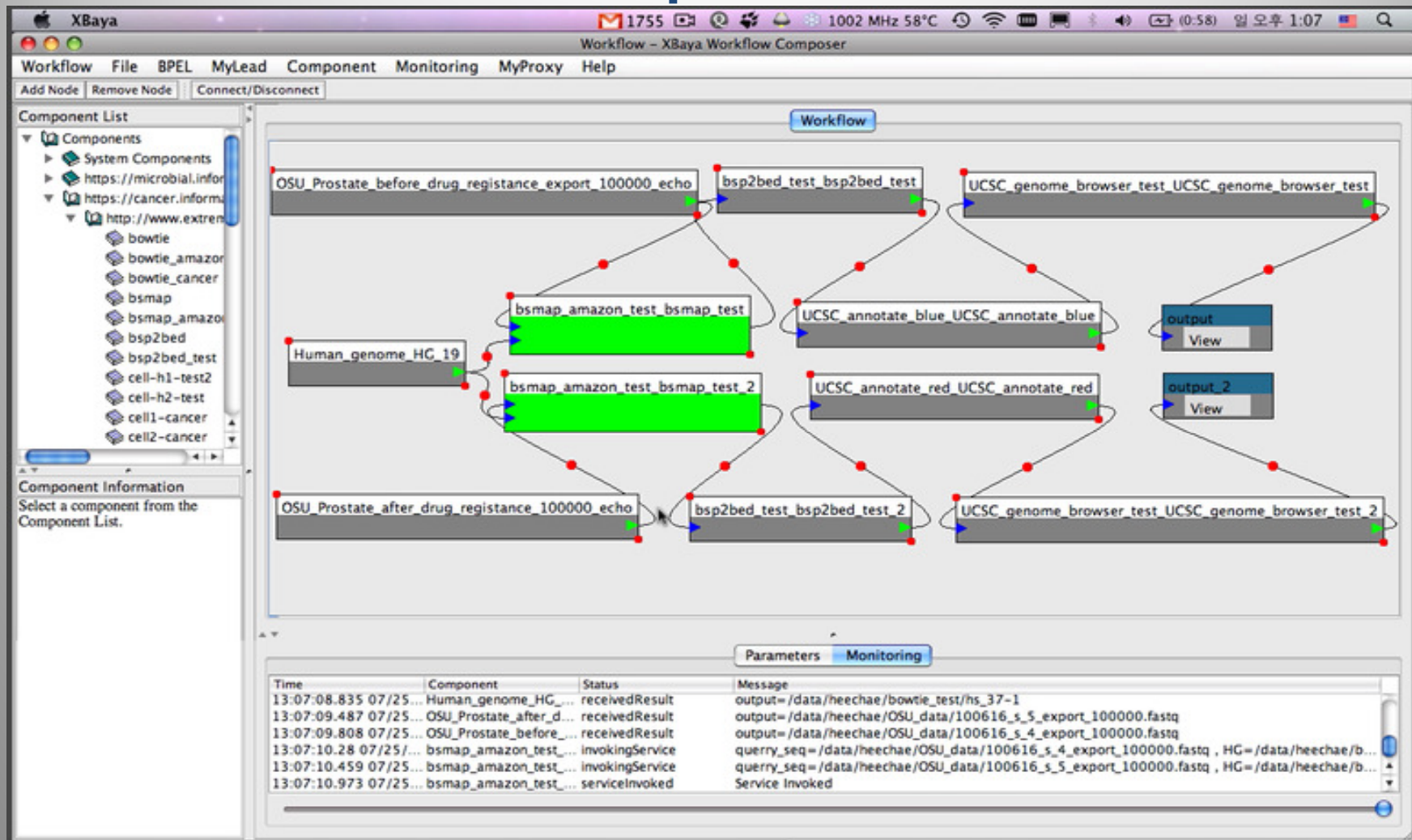


Data (in green color) is ready.

-- Sun Kim group at IU --



# BioVLAB-mCpG Screenshots

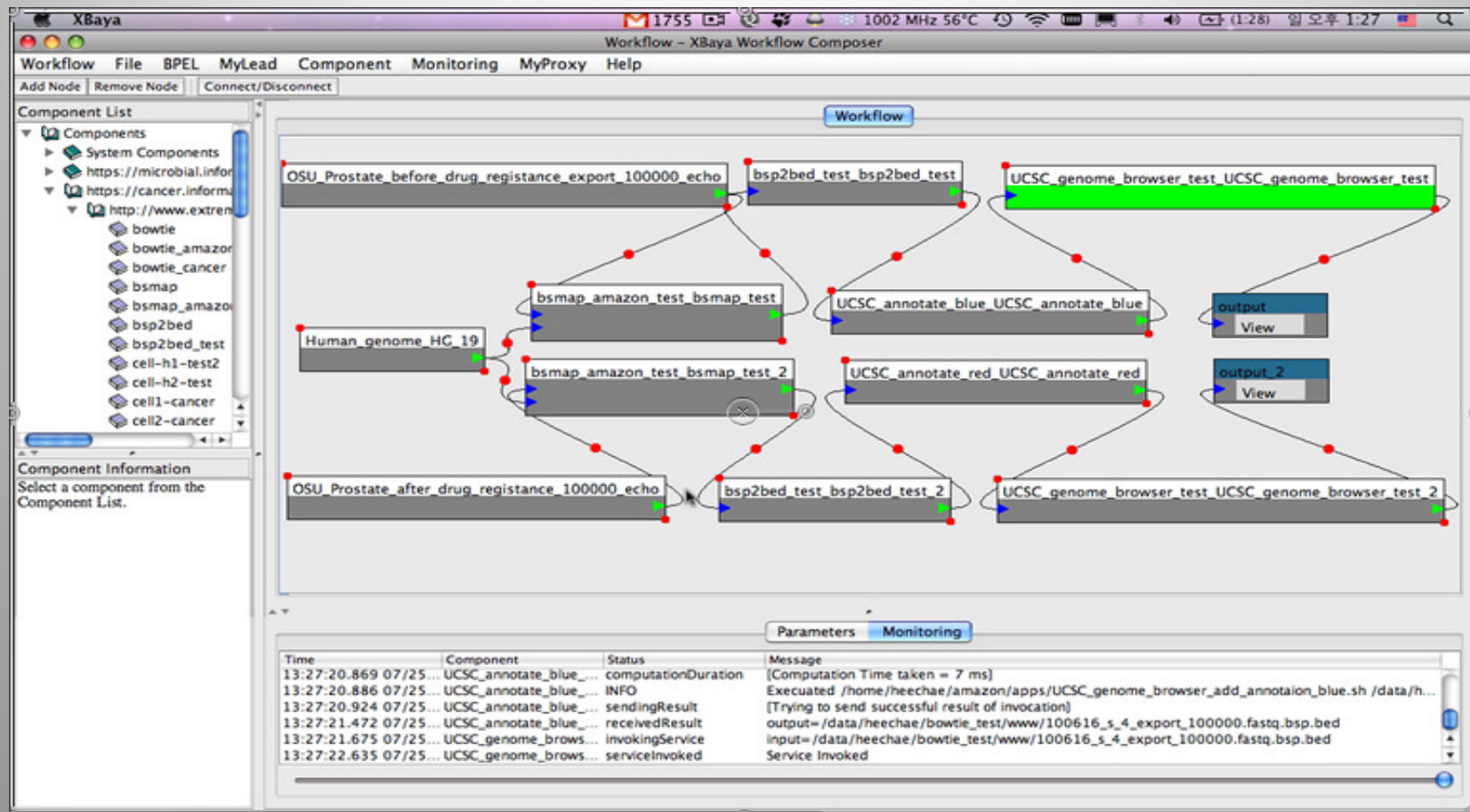


Sequence reads are being mapped by BSmap (green color).

-- Sun Kim group at IU --

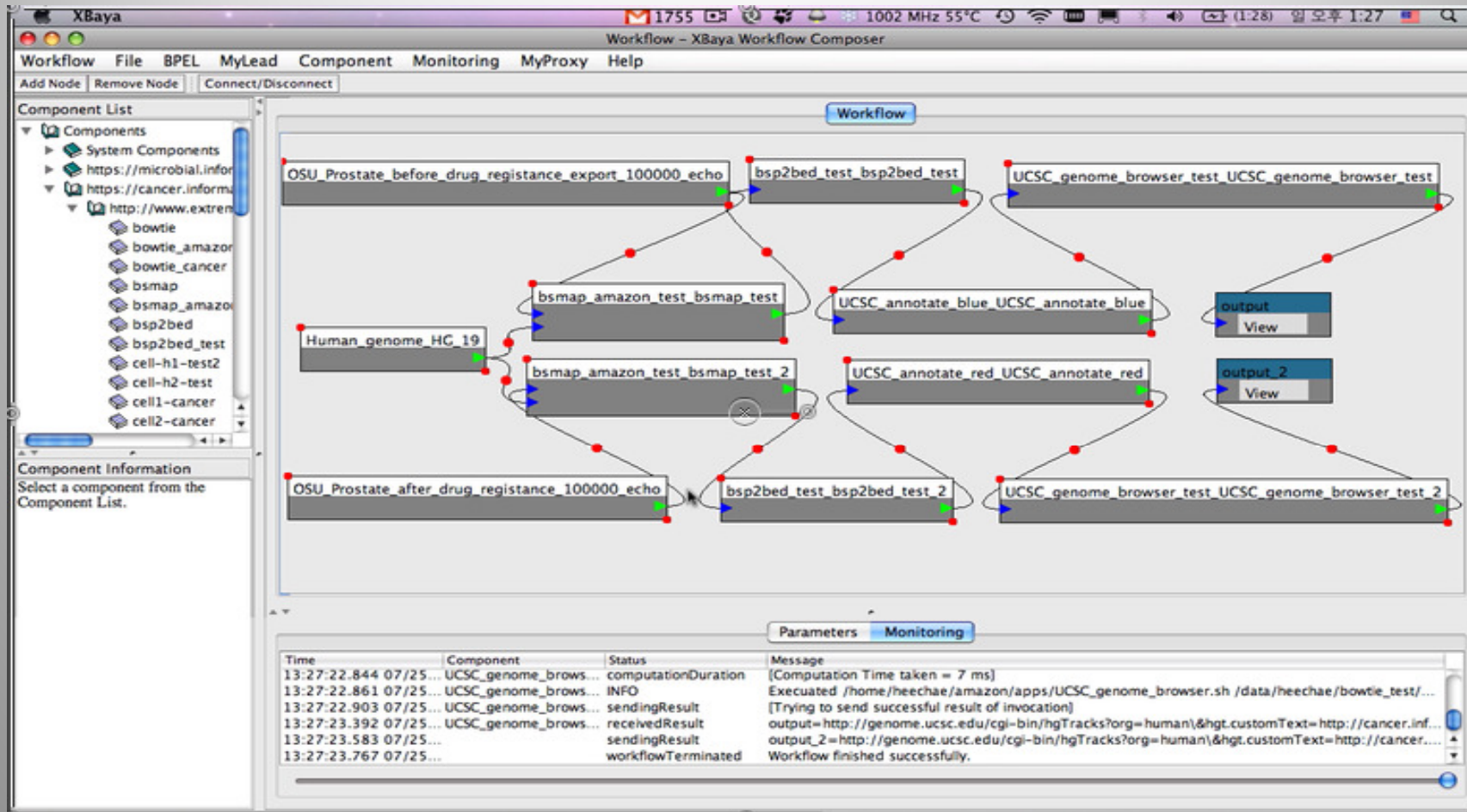


# BioVLAB-mCpG Screenshots



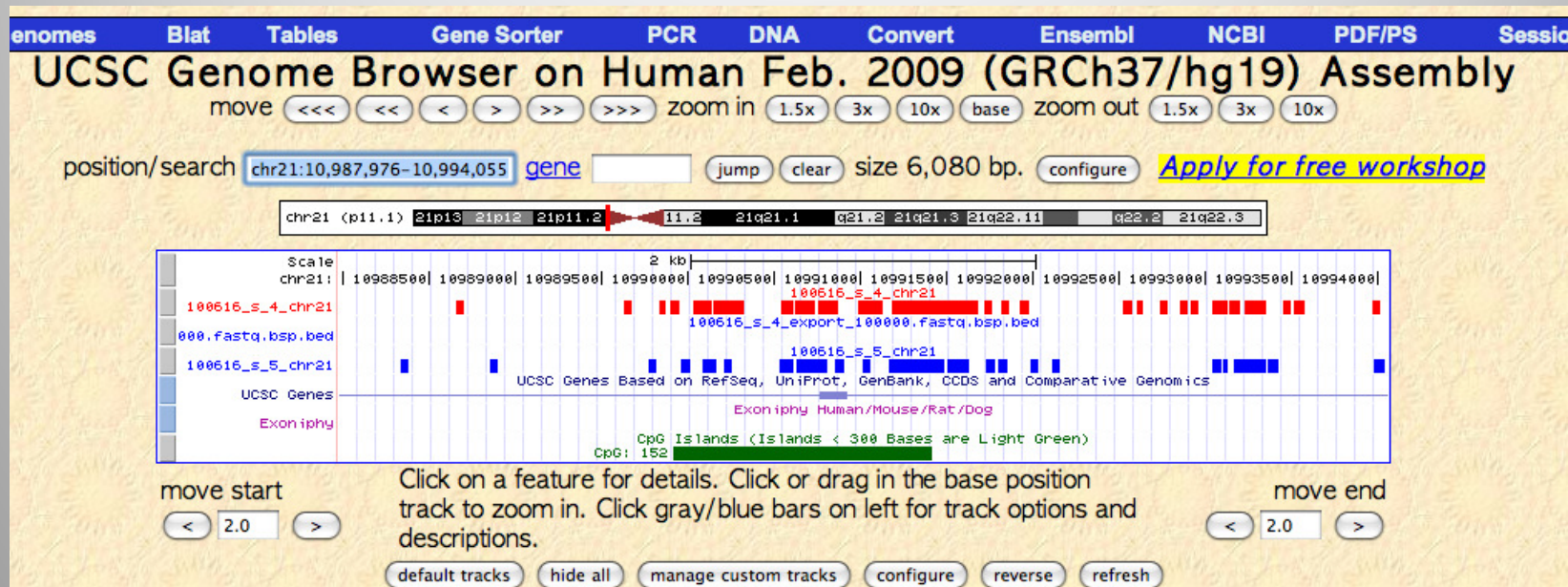
Uploading the result to the UCSC Genome Browser. (green color).

# BioVLAB-mCpG Screenshots



Finished! Let's look at visualized data.

# BioVLAB-mCpG Screenshots



Two lines (in red and blue colors) show DNA methylation status in the context of exon and a CpG Island.

# Acknowledgements

- Heejoon Chae, Youngik Yang, Hyungro Lee, Jong Yul Choi
- Suresh Marru, Chathura Herath, Marlon Pierce
- Ken Nephew at IU, Tim Huang at OSU and OSU-IU CCSB members
- NCI ICBP
- TeraGrid
- IU UITs

Thank you!!