Genomic data analysis with R

Macrogen Inc 김세환

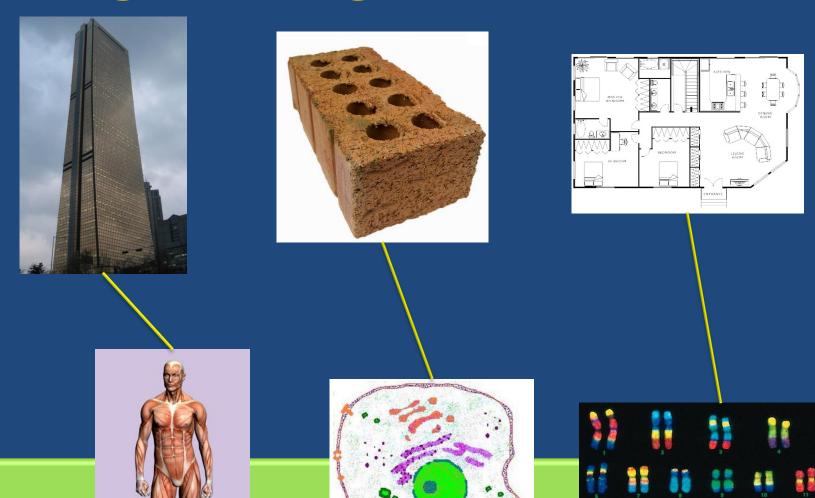
Contents

- Biological background
- Methods for biological data measurement
- Example of biological data analysis
- Bioconductor
- Trend in Genomics

Where We Are

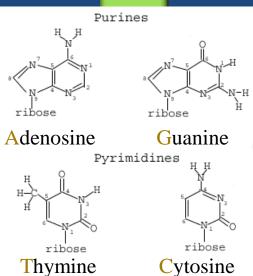
- Biological background
- Methods for biological data measurement
- Example of biological data analysis
- Bioconductor
- Trend in Genomics

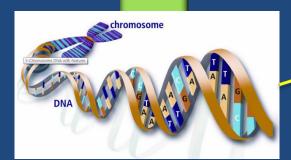
Biological Background for Genomics



Biological Background for Genomics









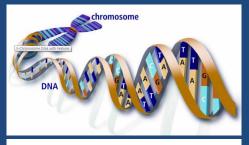
46 chromosomes = 22 * 2 + (1+1)

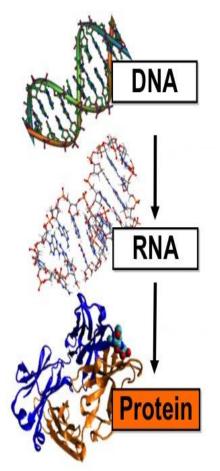


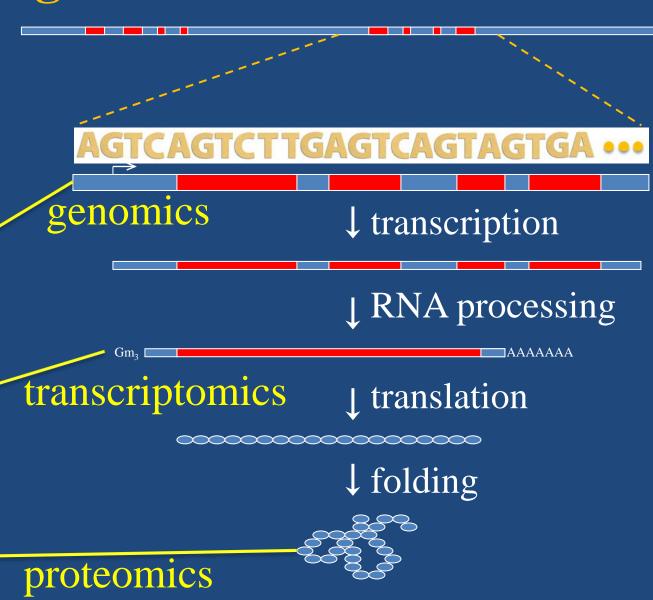
4 bases

 $3 * 10^9 \text{ bp}$

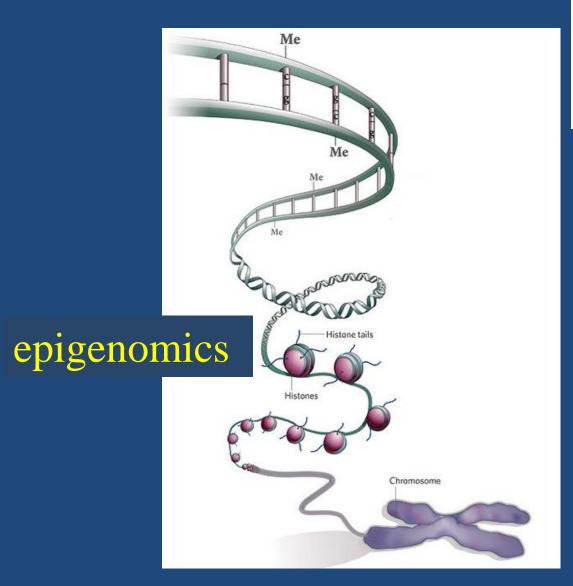
Central Dogma

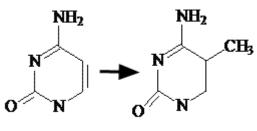


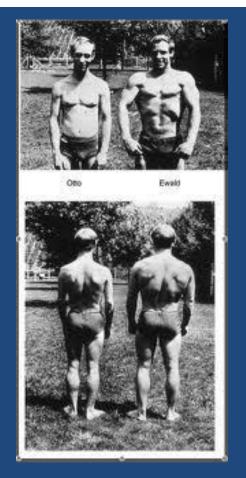




Methylated 5'C; Fifth base in DNA







Where We Are

- Biological background
- Methods for biological data measurement
- Example of biological data analysis
- Bioconductor
- Trend in Genomics

How to get these biological data?

genomics transcriptomics epigenomics

Measurement Methods

Measurement Methods

•Microarray

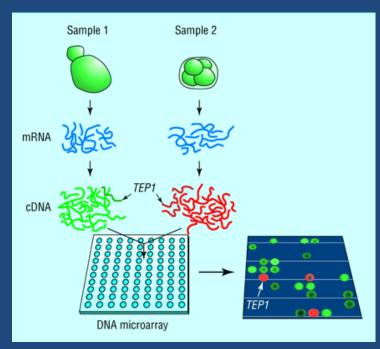


•DNA sequencer

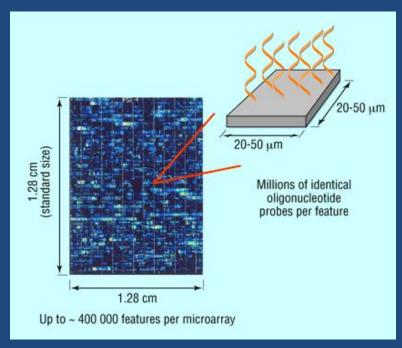


Microarray

- Probe => Measurement
- Analogue signal
- oligonucleotide array VS two-color spotted array





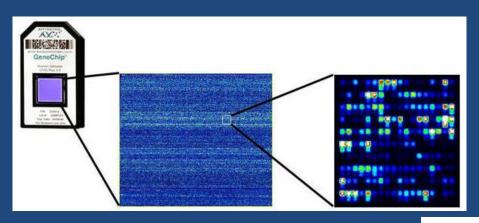


認用類的問題

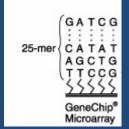
認。同樣的問

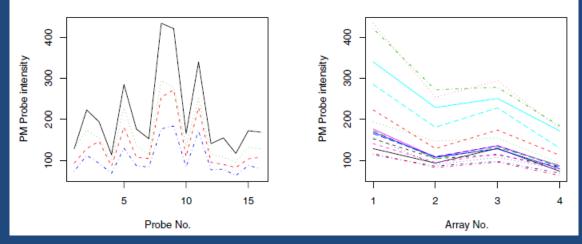
<oligonucleotide array>

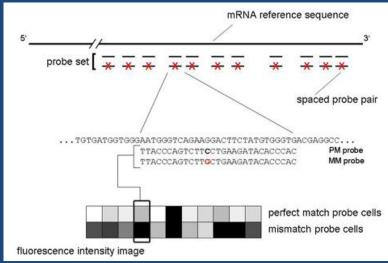
Preprocessing of Microarray Data

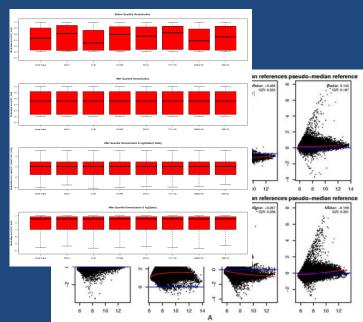


affy, affyPLM, affycomp, gcrma, affypdnn



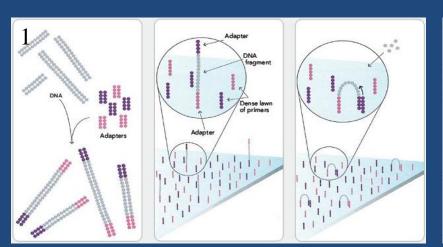


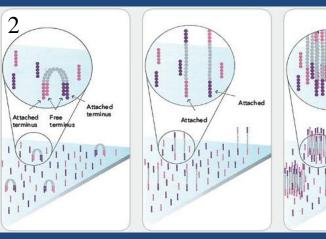


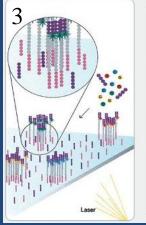


DNA sequencer

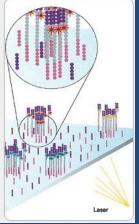
- •Measurement $=> \overline{\text{Probe}}$
- •Digital signal

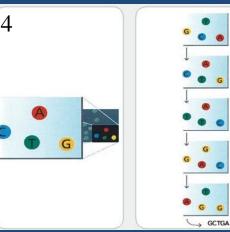


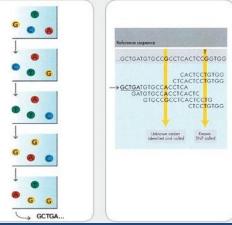














Preprocessing of Sequence Data



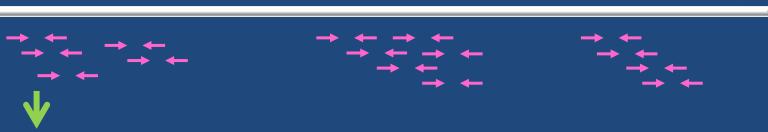
@HWUSI-EAS052R_0029:3:1:4952:1104#1/1 CTATCATGATCACCAACATCACCATCACC +HWUSI-EAS052R_0029:3:1:4952:1104#1/1 gggggfffefba_caB``abbbabaeeee

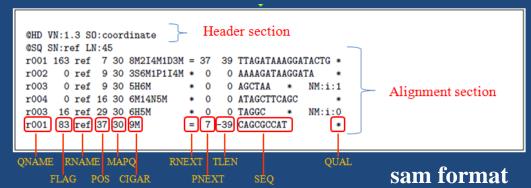


@HWUSI-EAS052R_0029:3:1:4952:1104#1/1 CTATCATGATCACCAACATCACCATCACCANGAACN +HWUSI-EAS052R_0029:3:1:4952:1104#1/1 gggggfffefba_caB``abbbabaeeeecB_\ZaB

Fastq file format

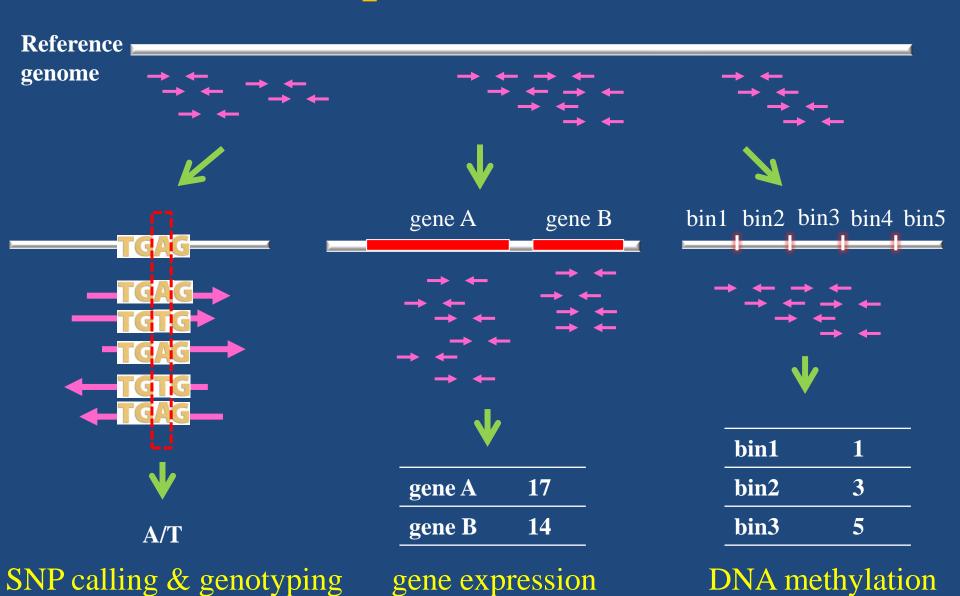






shortRead, Rsamtools

Data from Sequencer



Data Formats

- SNP (A, C, G, T): categorical
 - Genotype (AA, AB, BB)
 - Sample x Genotype for each SNP locus
 - Genetic Model(dominant, recessive, co-dominant, additive, allele models)

Expression

Signal to noise for each sample and probe (continuous)

• Methylation: categorical or continuous

- Avg_beta, methylation level for each CpG site(continuous)
- methylated/unmethylated (categorical) or hypo methylated, medium methylated, hyper methylated (categorical)

Where We Are

- Biological background
- Methods for biological data measurement
- Example of biological data analysis
- Bioconductor
- Trend in Genomics

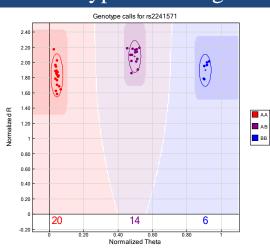
Genomics Data Analysis

	Normal				Patient	
Α	Α	Α	SNP1	Α	Α	Α
Т	Т	Т	SNP2	Т	Т	Т
G	G	G	SNP3	G	G	G
С	С	С	SNP4			
G	G	G	SNP5	G	G	G
G	G	G	SNP6	G	G	G
С	С	С	SNP7	С	С	С
	•				•	
•	•	•	•	•	•	•
C	C	С	SNP8	С	C	C
T	T	Т	SNP9	T	T	Т
T	Т	T	SNP10	T	Т	Т

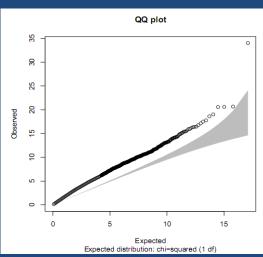
Genomics Data Analysis



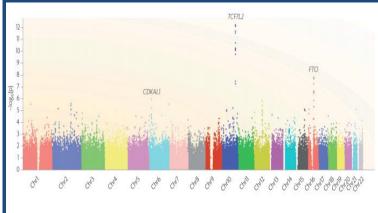
Genotype clustering



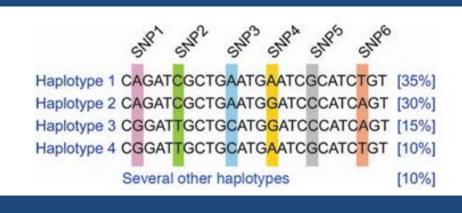
■ HWE test

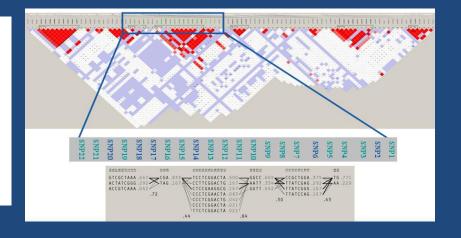


Manhattan Plot



LD & Haplotype



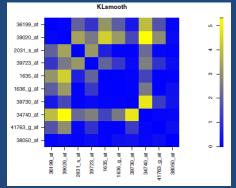


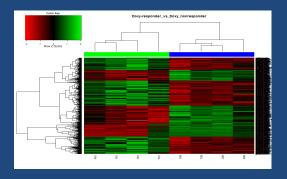
Transcriptomics/Epigenomic Data Analysis

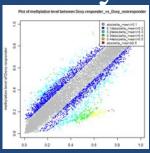
- Distance between Gene or Samples
 - Minkowski metric distance
 - Correlation-based distance
 - Distance between distributions
- Cluster analysis —> cluster, hopach
 - Partitioning
 - Hierachical
 - Hybrid
- Analysis of Differential Gene Expression/methylation

limma, multtest

- parametric test
- non-parametric test
- multiple test adjustment







Transcriptomics Data Analysis

normal



	Sample1	Sample2	Sample3
Gene A	13	14	15
Gene B			

Gene N

Non-specific filtering

	Sample 1	Sample 2	Sample 3
Gen e A	13	14	15
Gen e C			

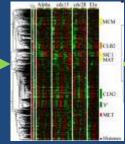
patient



	Sample1	Sample2	Sample3
Gene A	4	5	3
Gene B			

Gene N

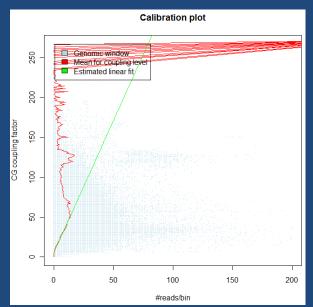
clustering

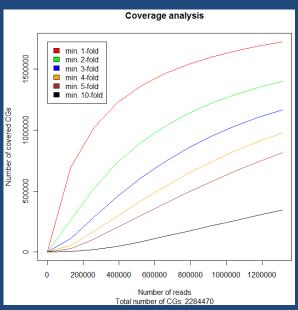


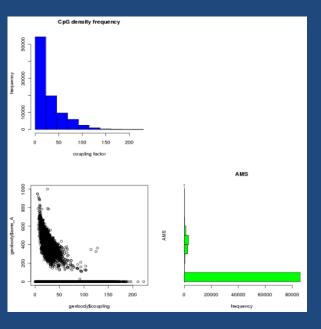
Finding DEG

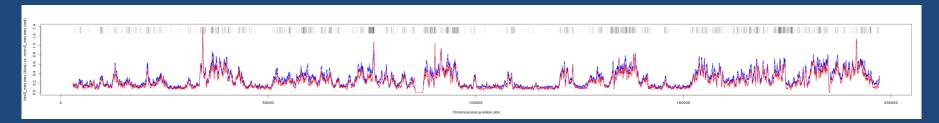
Gene K Gene I Gene M

Epigenomics Data Analysis



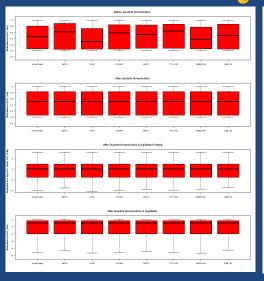


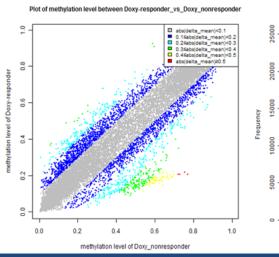


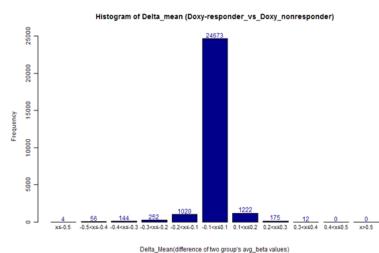


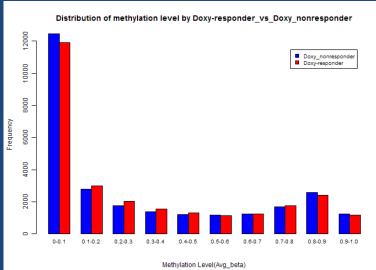


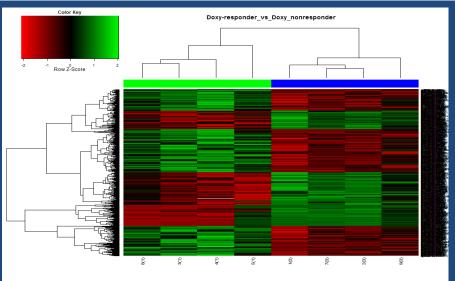
Transcriptomics/Epigenomic Data Analysis



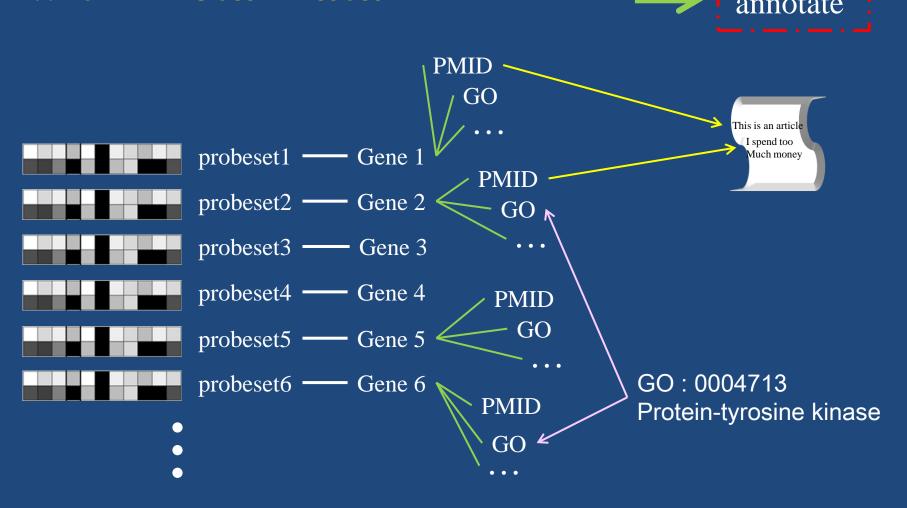








Association Genomic Data with Meta-Data



Where We Are

- Biological background
- Methods for biological data measurement
- Example of biological data analysis
- Bioconductor
- Trend in Genomics

Bioconductor

Method

Open Access

Bioconductor: open software development for computational biology and bioinformatics

Robert C Gentleman¹, Vincent J Carey², Douglas M Bates³, Ben Bolstad⁴, Marcel Dettling⁵, Sandrine Dudoit⁴, Byron Ellis⁶, Laurent Gautier⁷, Yongchao Ge⁸, Jeff Gentry¹, Kurt Hornik⁹, Torsten Hothorn¹⁰, Wolfgang Huber¹¹, Stefano Iacus¹², Rafael Irizarry¹³, Friedrich Leisch⁹, Cheng Li¹, Martin Maechler⁵, Anthony J Rossini¹⁴, Gunther Sawitzki¹⁵, Colin Smith¹⁶, Gordon Smyth¹⁷, Luke Tierney¹⁸, Jean YH Yang¹⁹ and Jianhua Zhang¹

Addresses: Department of Biostatistical Science, Dana-Farber Cancer Institute, 44 Binney St, Boston, MA 02115, USA. 2Channing Laboratory, Brigham and Women's Hospital, 75 Francis Street, Boston, MA 02115, USA. 3Department of Statistics, University of Wisconsin-Madison, 1210 W Dayton St, Madison, WI 53706, USA. 4Division of Biostatistics, University of California, Berkeley, 140 Warren Hall, Berkeley, CA 94720-7360. USA, 5Seminar for Statistics LEO C16. ETH Zentrum, Zürich CH-8092, Switzerl. 6Department of Statistics, Harvard University, 1 Oxford St, Cambridge, MA 02138, USA. 7Center for Biological Sequence Analysis, Technical University of Denmark, Building 208, Lyngby 2800, Denmark. 8Department of Biomathematical Sciences, Mount Sinai School of Medicine, 1 Gustave Levy Place, Box 1023, New York, NY 10029, USA. 9Institut für Statistik und Wahrscheinlichkeitstheorie, TU Wien, Wiedner Hauptstrasse 8-10/1071, Wien 1040, Austria. 10Institut für Medizininformatik, Biometrie und Epidemiologie, Friedrich-Alexander-Universität Erlangen-Nürnberg, Waldstraße6, D-91054 Erlangen, Germanv, 11Division of Molecular Genome Analysis, DKFZ (German Cancer Research Center), 69120 Heidelberg, Germanv, 12Department of Economics, University of Milan, 23 Via Mercalli, I-20123 Milan, Italy. 13 Department of Biostatistics, Johns Hopkins University, 615 N Wolfe St E3035, Baltimore, MD 21205, USA, 4Department of Medical Education and Biomedical Informatics, University of Washington, Box 357240, 1959 NE Pacific, Seattle, WA 98195, USA. 15Statistisches Labor, Institut für Angewandte Mathematik, Im Neuenheimer Feld 294, D 69120, Heidelberg, Germany. 16 Department of Molecular Biology, The Scripps Research Institute, 10550 North Torrey Pines Road, TPC-28, La Jolla, CA 92037, USA. 17Division of Genetics and Bioinformatics, The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville, Victoria 3050, Australia. 18Department of Statistics and Actuarial Science, University of Iowa, 241 Schaeffer Hall, Iowa City, IA 52242, USA. 49Center for Bioinformatics and Molecular Biostatistics, University of California, San Francisco, 500 Parnassus Ave, San Francisco 94143-0560, USA.

Correspondence: Robert C Gentleman. E-mail: rgentlem@jimmy.harvard.edu

Published: 15 September 2004

Genome Biology 2004, 5:R80

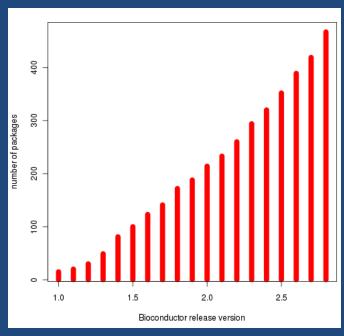
The electronic version of this article is the complete one and can be found online at http://genomebiology.com/2004/5/10/R80

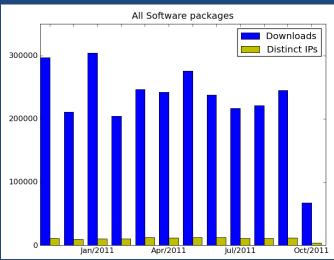
Received: 19 April 2004 Revised: 1 July 2004 Accepted: 3 August 2004 **Purpose**

Merit

expreSet class

Bioconductor





Bioconductor version 2.8 (Release) ▼ AnnotationData (593) ▶ ChipManufacturer (356) ▶ ChipName (190) CustomArray (2) ▶ CustomCDF (16) ▶ CustomDBSchema (10) FunctionalAnnotation (10) ▶ Organism (424) ▶ SequenceAnnotation (2) ▼ ExperimentData (82) ▶ Cancer (18) ChIPchipData (1) ChIPseqData (3) EColiData (1) FlowCytData (1) HapMap (7) HighThroughputSequencingData (1) HIV (1) MassSpectrometryData (1) NormalTissue (1) RNAExpressionData (1) RNAseqData (1) StemCells (1) Yeast (9) Software (467) ▶ Annotation (61) ▶ AssayDomains (182) ▶ AssayTechnologies (290) ▶ Bioinformatics (273) ▶ BiologicalDomains (46) ▶ Infrastructure (194)

Where We Are

- Biological background
- Methods for biological data measurement
- Example of biological data analysis
- Bioconductor
- Trend in Genomics

Why is genomics highlighted?

2001

Human Genome Project





2007

Next Generation Sequencer



1st generation sequencer







next generation sequencer

2011..

Genomics in Business

















Thank you

GATTGGGGTTTTCCCCTCCCATGTGCTCAAGACTGGCGCTAAAAGTTTTGAGCTTCTCAAAAGTCTAGAG CCACCGTCCAGGGAGCAGGTAGCTGCTGGGCTCCGGGGACACTTTGCGTTCGGGCTGGGAGCGTGCTT TCCACGACGGTGACACGCTTCCCTGGATTGGgtaagctcctgactgaacttgatgagtcctctctgagtcacgggctctcggt ataggagttgcattgttgggagacctgggtgtagatgatggggatgttaggaccatccgaactcaaagttgaacgcctaggcagaggagtggagctttggg gaaccttgagccggcctaaagcgtacttctttgcacatccacccggtgctgggcgtagggaatccctgaaataaaagatgcacaaagcattgaggtctga gacttttggatctcgaaacattgagaactcatagctgtatattttagagcccatggcatcctagtgaaaactggggctccattccgaaatgatcatttgggggtg tccaattgaaggctgtcagtcgtggaagtgagaagtgctaaaccaggggtttgcccgccaggccgaggaggaccgtcgcaatctgagaggcccggcag ggtacgtctgagaatcaaattttgaaagagtgcaatgatgggtcgtttgata<mark>att</mark>tgtcggaaaaacaat**cta**cctgttatc**tag**ctttggg**cta**ggccattcca ggtaggaggcggaactcgaattcatttctcccgctgccccatctcttagctcgcggttgtttcattccgcagtttcttcccatgcacctgccgcgtaccggccact ttgtgccgtacttacgtcatctttttcctaaatcgaggtggcatttacacacagcgcc<mark>agt</mark>gcacacag**caagtgc**acagga**agatgagtttt**ggcccctaac cgctccgtgatgcctaccaagtcacagacccttttcatcgtcccagaaacgtttcatcacgtctcttcccagtcgattcccgaccccacctttattttgatctccat gataggagttccagaccagcgtggccaacgtggtgaatccccgtctctactaaaaaatacaaaaattagctgggcgtggtggtggctgtaatcccagcta ttcgggagggtgaggcaggagaatcgcttgaacccgggaggcagaggttgcagtgagccaagatcgtgccactacactccagcctgggcgacaaga acgaaactccgtctcaaaaaaaaggggggaatcatacattatgtgctcatttttgtcgggcttctgtccttcaatgtactgtctgacattcgttcatgttgtatatat cagtattttgctccttttcatttagtatagtccatcgattgtatatccgtccttttgatggccttttgagttgtttcccatttgcggttatgaaataaagctgctataaacatt cttgtacaattctttttgtgatcatatgttttcgtgtttcttggagaaatacttaggagggaattgcgagtttggaagtaaaaagtagctgtattttgaactttttcaga