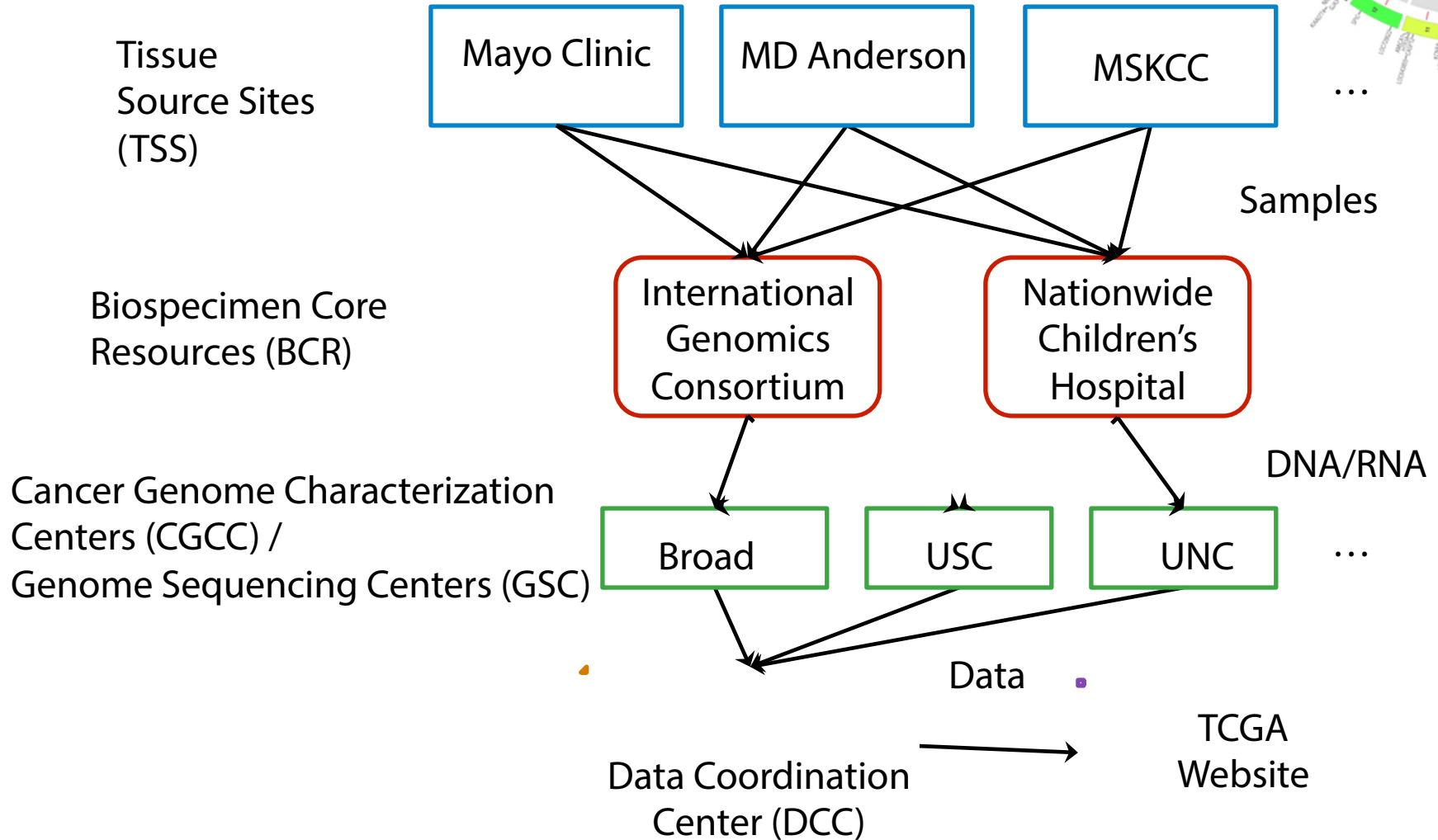The Cancer Genome Atlas

# Detection, Diagnosis <u>and Correction</u> of Batch Effects in TCGA Data

*Rehan Akbani*

*Dept. of Bioinformatics and Computational Biology,*
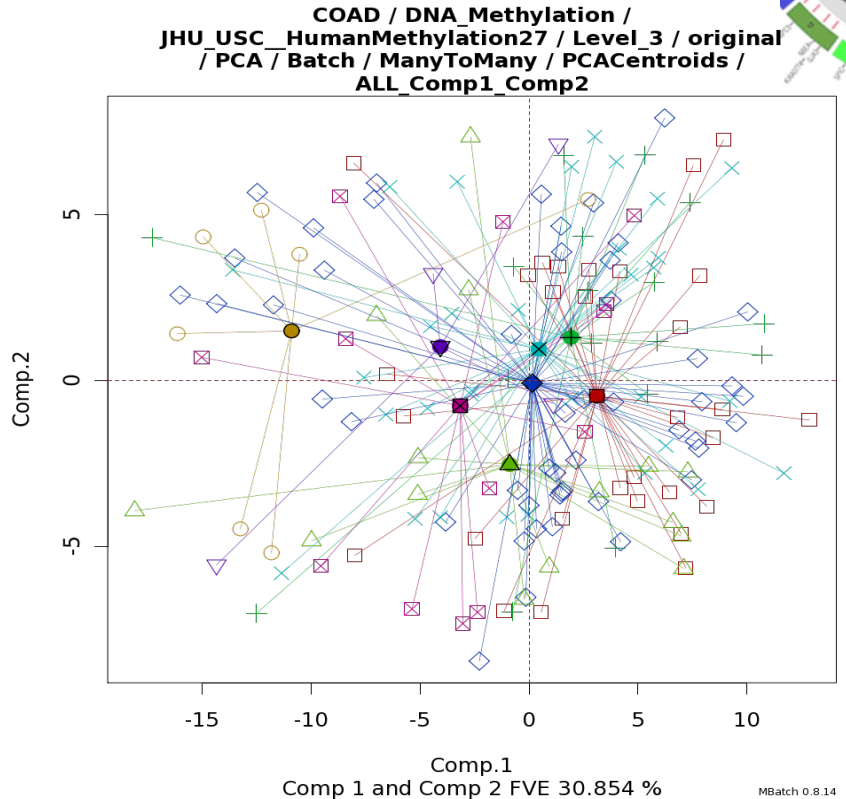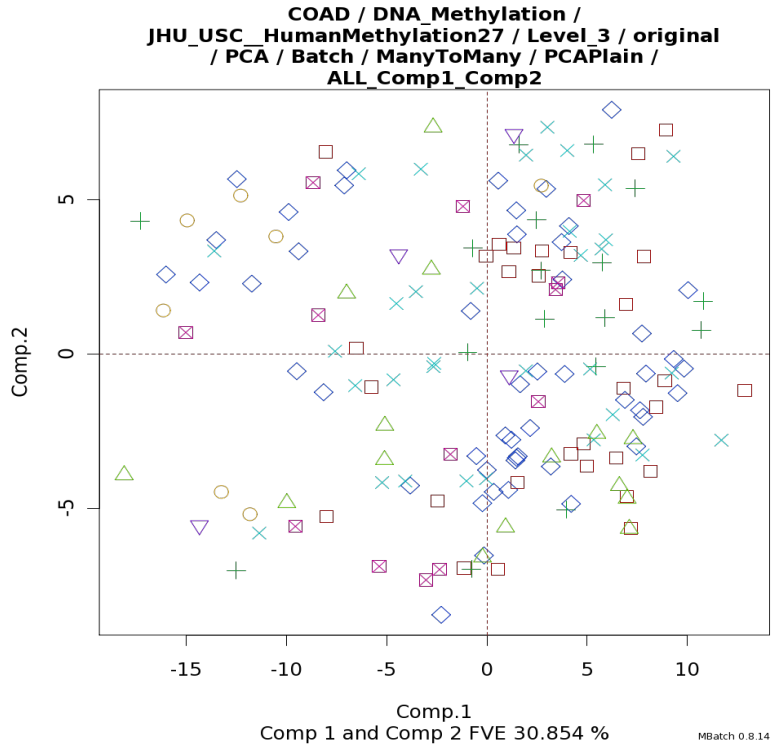*UT MD Anderson Cancer Center*

THE UNIVERSITY OF TEXAS
MD Anderson
Cancer Center

# Simplified Flow Diagram for TCGA



Tissue Source Sites (TSS)

Mayo Clinic    MD Anderson    MSKCC    …

Samples

Biospecimen Core Resources (BCR)

International Genomics Consortium    Nationwide Children's Hospital

Cancer Genome Characterization Centers (CGCC) / Genome Sequencing Centers (GSC)

DNA/RNA

Broad    USC    UNC    …

Data

Data Coordination Center (DCC)

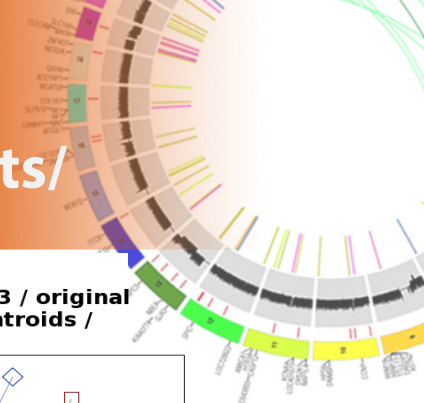TCGA Website

The Cancer Genome Atlas

# Step 1: Batch effects diagnoses

- Objectives
  - Detect / quantify batch effects
  - Identify source(s) of batch effects
- Tools / algorithms for diagnoses – MBatch R package

  http://bioinformatics.mdanderson.org/tcgabatcheffects/

  1. PCA-Plus plots (novel)
  2. BatchCorr algorithm (novel)
  3. Hierarchical clustering
  4. Clinical correlates
  5. Box plots
  6. ANOVA / MANOVA

- Disclaimer: No substitute for human input

The Cancer Genome Atlas

# 1. PCA-Plus

## http://bioinformatics.mdanderson.org/tcgabatcheffects/



COAD / DNA_Methylation /
JHU_USC__HumanMethylation27 / Level_3 / original
/ PCA / Batch / ManyToMany / PCAPlain /
ALL_Comp1_Comp2

Comp.1
Comp 1 and Comp 2 FVE 30.854 %

MBatch 0.8.14

COAD / DNA_Methylation /
JHU_USC__HumanMethylation27 / Level_3 / original
/ PCA / Batch / ManyToMany / PCACentroids /
ALL_Comp1_Comp2

Comp.1
Comp 1 and Comp 2 FVE 30.854 %

MBatch 0.8.14

### Batch (Points)

- 28 (30)
- 29 (7)
- 30 (15)
- 33 (17)
- 36 (32)
- 41 (49)
- 45 (4)
- 66 (13)

### Dispersion Metrics

First PCA Component: 1
First Component FVE (%): 22.972
Second PCA Component: 2
Second Component FVE (%): 7.882

DSC(1,2): 0.411
Dw(1,2): 7.51
Db(1,2): 3.087
DSC pvalue(1,2): 0

DSC: 0.299
Dw: 14.003
Db: 4.193
DSC pvalue: 0

The Cancer Genome Atlas

4

# 2. BatchCorr



BatchCorr value

| | 28 | 29 | 30 | 33 | 36 | 41 | 45 | 66 |
|---|---|---|---|---|---|---|---|---|
| CZN | 0.75 | 0.383 | 0.729 | 0.826 | 0.957 | 0.877 | 0.63 | 0.941 |
| CZR | 0.746 | 0.446 | 0.928 | 0.871 | 1 | 0.863 | 0.314 | 1 |

p-values

| | 28 | 29 | 30 | 33 | 36 | 41 | 45 | 66 |
|---|---|---|---|---|---|---|---|---|
| CZN | 0.048 | 0.178 | 0.34 | 0.59 | 0.867 | 0.021 | 0.606 | 0.967 |
| CZR | 0.045 | 0.195 | 0.657 | 0.46 | 1 | 0.016 | 0.818 | 1 |

Batch effects present:

BatchCorr < 0.7 AND
p-value < 0.05

The Cancer Genome Atlas

# 3. Hierarchical Clustering



COAD/Expression-Genes/UNC__AgilentG4502A_07/Level_3

4. Clinical Correlates COAD/READ miRNA

# 5. Box plots



Batch medians

Batch means

The Cancer Genome Atlas

The Cancer Genome Atlas

# Step 2: Batch effects correction

- Correct the source of the problem whenever possible

- When not possible, or source unknown, algorithms can be used

- Some algorithms for correction:
  - ComBat (aka Empirical Bayes)
  - ANOVA
  - Median Polish

- Included in MBatch R package

  http://bioinformatics.mdanderson.org/tcgabatcheffects/

The Cancer Genome Atlas

# Kidney cancer (KIRC) DNA methylation data (27k)

# After removing sex chromosomes



New dichotomy based on batch ID appears

Legend:
- 32 (blue, ○)
- 50 (red, △)
- 64 (green, +)
- 65 (purple, ×)
- 69 (yellow, ◇)

After removing Sex Chr by Batch

No Sing-Cha, By Batch

Second Comp.

First Comp.

Chromophobes

32
50
64
65
69

# TCGA MBatch website
## http://bioinformatics.mdanderson.org/tcgabatcheffects/

# TCGA MBatch website
## http://bioinformatics.mdanderson.org/tcgabatcheffects/

# TCGA MBatch website
## http://bioinformatics.mdanderson.org/tcgabatcheffects/

# TCGA MBatch website
## http://bioinformatics.mdanderson.org/tcgabatcheffects/

# TCGA MBatch website
## http://bioinformatics.mdanderson.org/tcgabatcheffects/

# Acknowledgments

UT MD Anderson

- Nianxiang Zhang
- Tod D. Casasent
- Chris Wakefield
- James M. Melott
- Anna K. Unruh
- Thomas C. Motter
- Bradley M. Broom
- John N. Weinstein

In-Silico

- James Cleland
- Andy Wong
- Mike Ryan

Poster # 50

The Cancer Genome Atlas