

Some challenges related to next-generation sequencing data

Rafael A. Irizarry

Department of Biostatistics

Johns Hopkins Bloomberg School of Public Health

<http://rafallab.org>

@rafallab

<http://genomics.jhu.edu>

Topics

- Base calling
- Biological variation and systematic biases
- Batch effects (Time permitting)
- Bump hunting (offline)

Base calling most slides courtesy of

Héctor Corrada Bravo and Joyce Hsiao
Center for Bioinformatics and Computational Biology
Dept. of Computer Science
University of Maryland-College Park

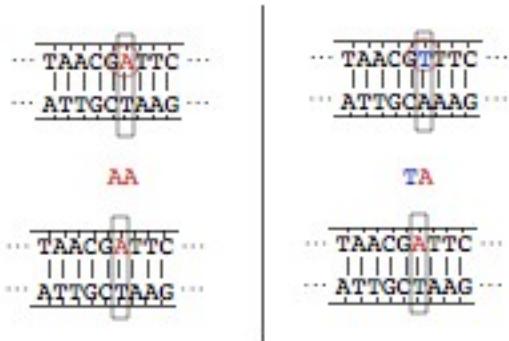
<http://cbcb.umd.edu/~hcorrada>



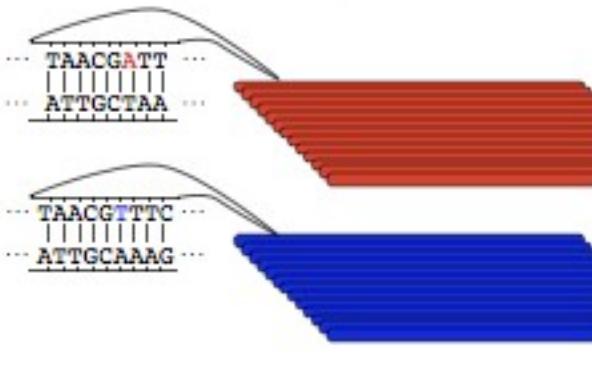


What is the basis of phenotypic variation?

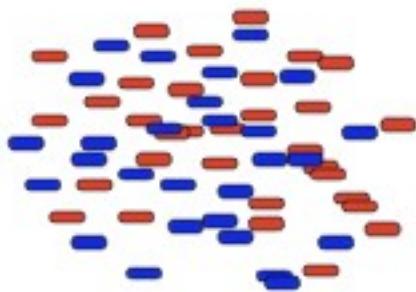
Genotyping



Sec-gen Sequencing for SNPs



Sec-gen Sequencing for SNPs



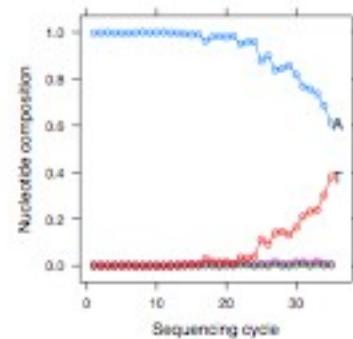
Sec-gen Sequencing for SNPs



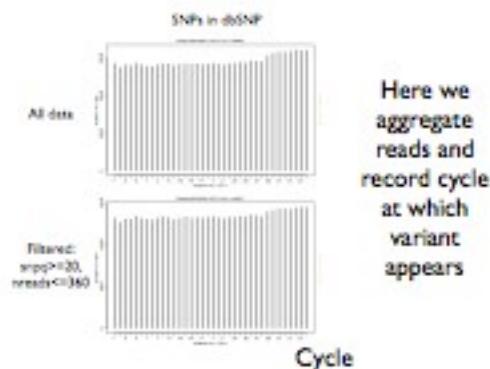
SNPs

SNPs

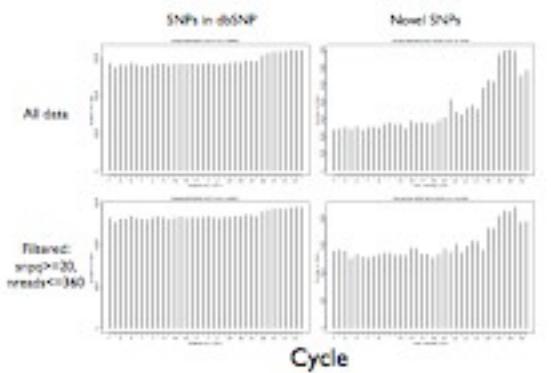
All Reads



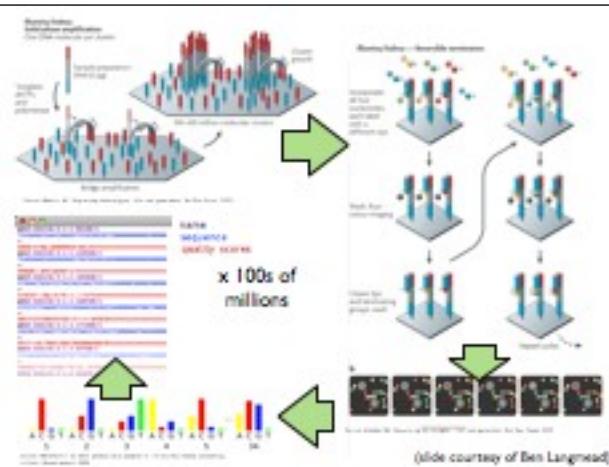
1000 Genomes Data



1000 Genomes Data



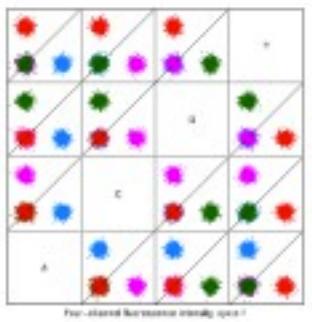
What is causing this?



Before Reads There were Intensities

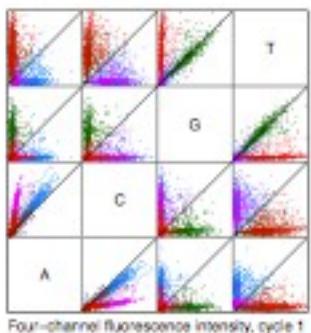
```
- 1002(1,18,C,A)  
A,1 C,1 G,1 T,1  
1 254.0 122.1 128.3 1380.9  
2 3953.5 3386.8 -767.4 300.3  
3 890.1 102.1 -927.4 101.1  
4 100.1 2007.9 -92.8 388.7  
5 979.4 6443.8 943.5 494.9  
6 540.1 8963.1 18.7 -1176.8  
7 253.0 253.0 25.5 4388.8  
8 3897.7 5894.5 -384.7 -94.1  
9 367.4 348.3 3886.2 5788.6  
10 1332.1 6424.4 -697.6 -349.2
```

We Want to See This



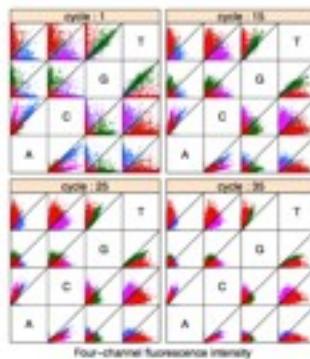
Color coded by call
made: A, C, G, T

But See This

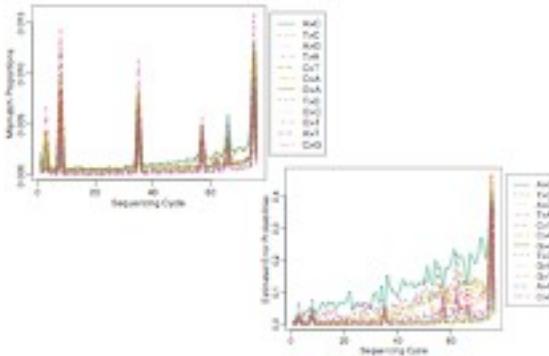


Color coded by call
made: A, C, G, T

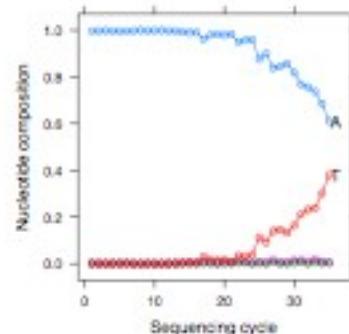
Gets Worse for higher cycles



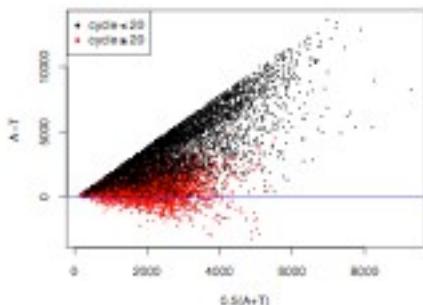
Error Rate and Reported Quality



Remember This!



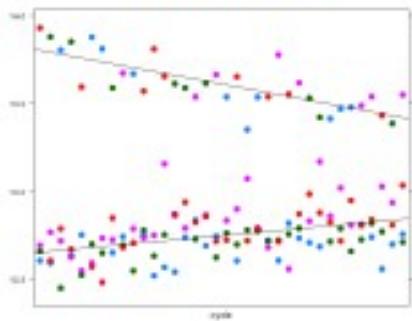
Bias Explained



Base Calling

- 1) Rougemont et al. Probabilistic base calling of Solexa sequencing data. *BMC Bioinformatics* (2008)
- 2) Ulrich et al. Alta-Cyclic: a self-optimizing base caller for next-generation sequencing. *Nat Methods* (2008)
- 3) Kao et al. BayesCall: A model-based base-calling algorithm for high-throughput short-read sequencing. *Genome Res* (2009)
- 4) Corrao Bracco and Istraty. Model-Based Quality Assessment and Base-Calling for Second-Generation Sequencing Data. *Biometrics* (2009)
- 5) Cokus et al. Shotgun bisulphite sequencing of the Arabidopsis genome reveals DNA methylation patterning. *Nature* (2009)

Intensity Model



Intensity Model

log intensity read i , cycle j , channel c

$$u_{ijc} = \underline{\Delta_{ijc}}(\mu_{cja} + \underline{x_j^T \alpha_i} + \underline{\epsilon_{ijc}^\alpha}) + \\ (1 - \underline{\Delta_{ijc}})(\mu_{cjb} + \underline{x_j^T \beta_i} + \underline{\epsilon_{ijc}^\beta})$$

indicators of nucleotide identity, read i , pos. j

$$\Delta_{ijc} = \begin{cases} 1 & \text{if } c \text{ is the nucleotide in read } i \text{ position } j \\ 0 & \text{otherwise} \end{cases}$$

Intensity Model

log intensity read i , cycle j , channel c

$$u_{ijc} = \underline{\Delta_{ijc}}(\mu_{cja} + \underline{x_j^T \alpha_i} + \underline{\epsilon_{ijc}^\alpha}) + \\ (1 - \underline{\Delta_{ijc}})(\mu_{cjb} + \underline{x_j^T \beta_i} + \underline{\epsilon_{ijc}^\beta})$$

read-specific linear models

Intensity Model

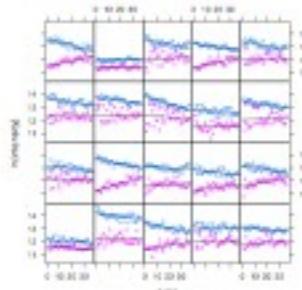
log intensity read i , cycle j , channel c

$$u_{ijc} = \underline{\Delta_{ijc}}(\mu_{cja} + \underline{x_j^T \alpha_i} + \underline{\epsilon_{ijc}^\alpha}) + \\ (1 - \underline{\Delta_{ijc}})(\mu_{cjb} + \underline{x_j^T \beta_i} + \underline{\epsilon_{ijc}^\beta})$$

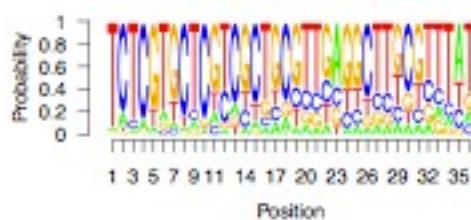
measurement error

$$\epsilon_{ijc}^\alpha \sim N(0, \sigma_{\alpha i}^2) \quad \epsilon_{ijc}^\beta \sim N(0, \sigma_{\beta i}^2)$$

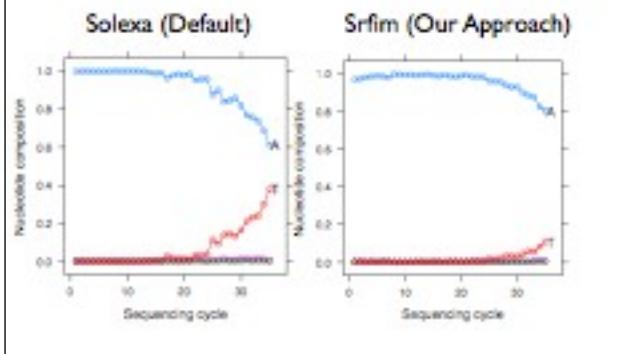
Read & Cycle Effects



Base Identity Probability Profiles



Before And After



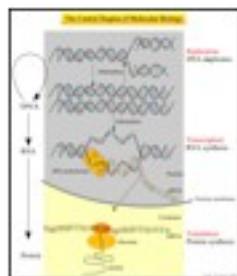
Biological variation and systematic biases

and its importance in gene expression
Examples from RNAseq data analysis

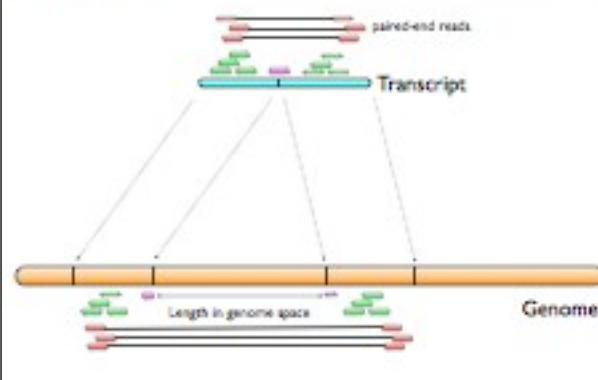
What is the basis of phenotypic variation?



Central Dogma of Biology



Mapping transcripts

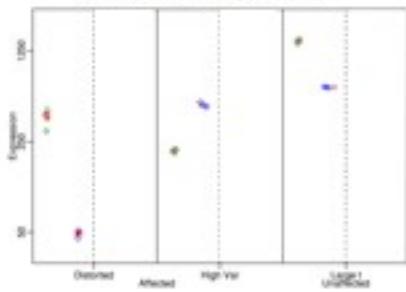


Most RNAseq studies do not include biological replicates

Hansen et al. (2011) Nature Biotechnology

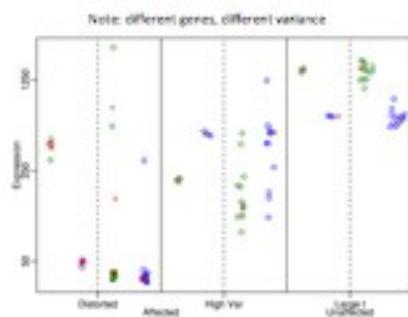
Biological versus technical replicates

Note: different genes, different variance

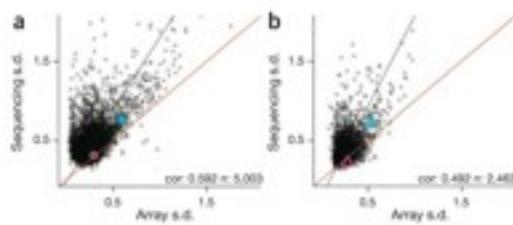


Kondratenko et al. (2008) PNAS

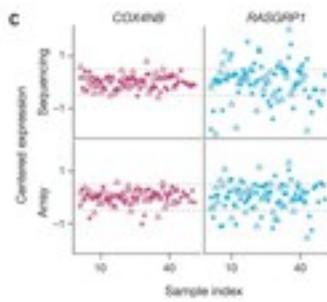
Biological versus technical replicates



SD NGS versus microarray



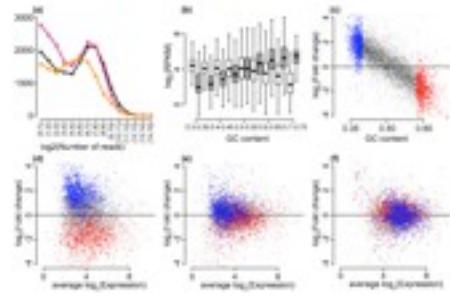
Highlighted genes by sample



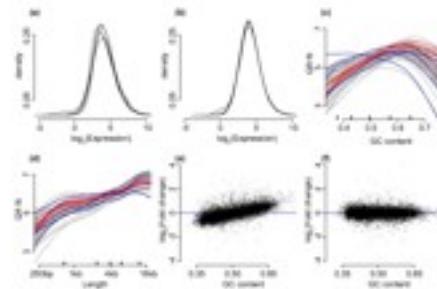
Systematic errors

Plots courtesy of Kasper Hansen and
Zhijin Wu

Different distributions and GC content effects



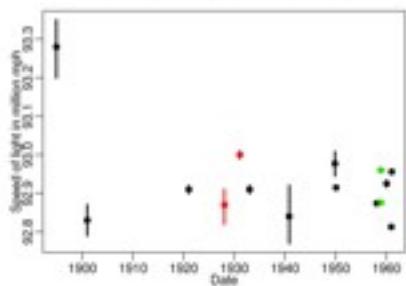
Estimated distributions and GC effects



Batch Effects

11

Speed of Light Estimates with "Confidence Intervals" (1900-1960)



Youden, WJ (1972). "Enduring Values." *Technometrics*, 14, 5-11.

Batch Effect in Genomics

Many slides provided by Jeff Leek

Gene Expression

LETTER

Nature Genetics 38, 402 - 403 (2006)
Published online 2 February 2006 | doi:10.1038/ng1766
Natural variation in human gene expression assessed in lymphoblastoid cells

Gene Expression

LETTER

Nature Genetics 38, 402 - 403 (2006)
Published online 2 February 2006 | doi:10.1038/ng1766
Natural variation in human gene expression assessed in lymphoblastoid cells

Correspondence

Nature Genetics 38, 807 - 808 (2006)
doi:10.1038/ng1707-807

On the design and analysis of gene expression studies in human populations

Joshua H Alvey¹, Shameek Basu², Jeffrey T Leek² & John D Storey^{1,2}

Proteomics

THE LANCET

3

© 2006 Elsevier Ltd. All rights reserved. 0950-2222/\$ - see front matter

doi:10.1016/j.lane.2006.01.001

Journal homepage: www.lancet.com

Letters - Monitoring of disease
Use of proteomic patterns in serum to identify ovarian cancer

Proteomics

THE LANCET
Volume 363, Issue 9399, 10 February 2009, Pages 503-508

Editorial - Mammalian Disease
Use of proteomic patterns in serum to identify ovarian cancer

See Single Protein in Cancer Testing (Editorial)

Photo

Science

Science is a weekly journal of original research, global news, and commentaries.

Author: ... et al.

Published online July 1, 2009
DOI: 10.1126/science.1160332

REPORT

Genetic Signatures of Exceptional Longevity in Humans

Paula J. Salomão^{1,2}, Maria Salomão², Anderson Puccio¹, Raphael M. Monteiro¹, Ethylmaia Matos¹, Silvia Andrade¹, Daniel S. Oliveira¹, Jozane R. Rezende¹, Renata H. Rezende¹, Monica M. Thomé¹, Henrique Thomé¹, Cláudia J. Reisendoerfer¹, and Thomas S. Perls^{2,3}

¹ Author Affiliations

* To whom correspondence should be addressed. E-mail: salomaom@uol.com.br (P.J.S.).

GWAS

Science

AAAS SCIENCE JOURNALS | SCIENCE & COMMUNITIES | MULTIMEDIA | COLLECTIONS

Science: The World's Leading Journal of Original Research, Global News, and Commentaries

Author: ... et al.

Published online July 1, 2009
DOI: 10.1126/science.1160332

REPORT

Genetic Signatures of Exceptional Longevity in Humans

Paula J. Salomão^{1,2}, Maria Salomão², Anderson Puccio¹, Raphael M. Monteiro¹, Ethylmaia Matos¹, Silvia Andrade¹, Daniel S. Oliveira¹, Jozane R. Rezende¹, Renata H. Rezende¹, Monica M. Thomé¹, Henrique Thomé¹, Cláudia J. Reisendoerfer¹, and Thomas S. Perls^{2,3}

¹ Author Affiliations

* To whom correspondence should be addressed. E-mail: salomaom@uol.com.br (P.J.S.).

GWAS

Science

AAAS SCIENCE JOURNALS | SCIENCE & COMMUNITIES | MULTIMEDIA | COLLECTIONS

Science: The World's Leading Journal of Original Research, Global News, and Commentaries

Author: ... et al.

This article has been retracted.

Published online July 1, 2009
DOI: 10.1126/science.1160332

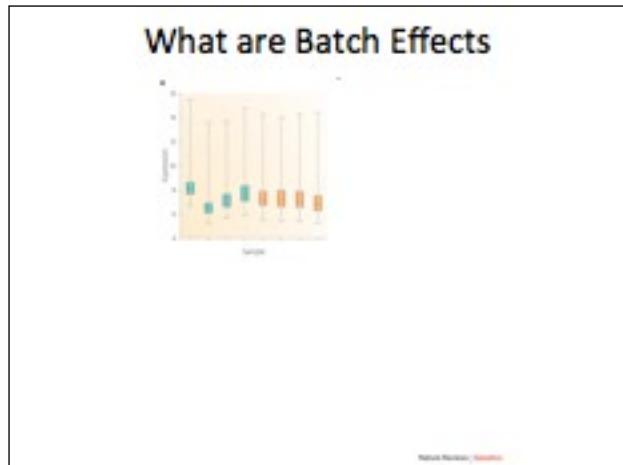
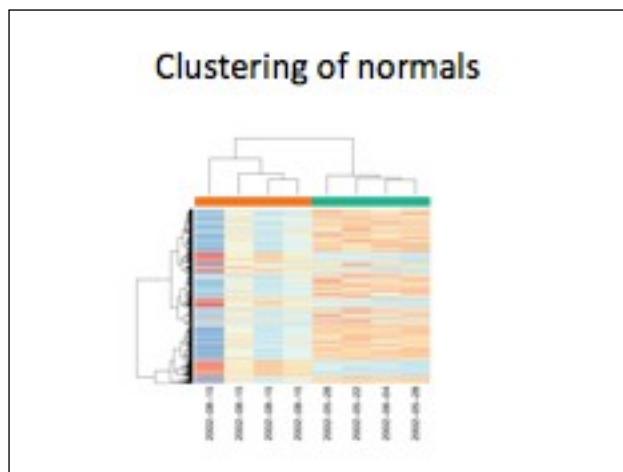
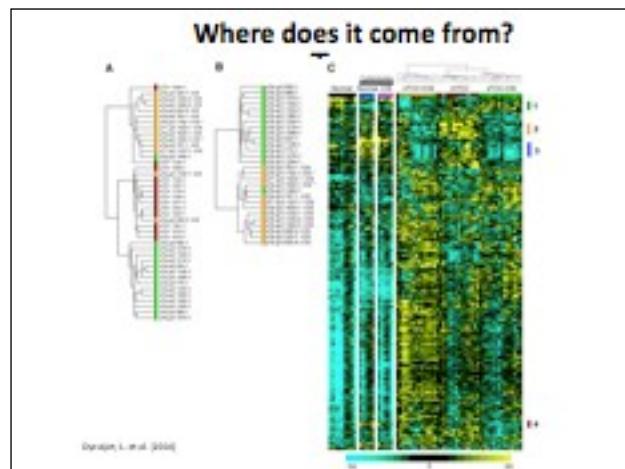
REPORT

Genetic Signatures of Exceptional Longevity in Humans

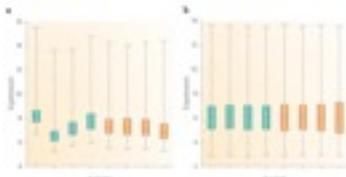
Paula J. Salomão^{1,2}, Maria Salomão², Anderson Puccio¹, Raphael M. Monteiro¹, Ethylmaia Matos¹, Silvia Andrade¹, Daniel S. Oliveira¹, Jozane R. Rezende¹, Renata H. Rezende¹, Monica M. Thomé¹, Henrique Thomé¹, Cláudia J. Reisendoerfer¹, and Thomas S. Perls^{2,3}

¹ Author Affiliations

* To whom correspondence should be addressed. E-mail: salomaom@uol.com.br (P.J.S.).

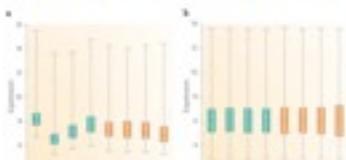


What are Batch Effects



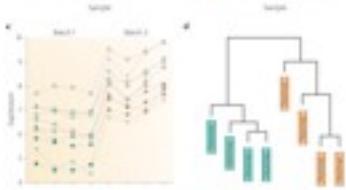
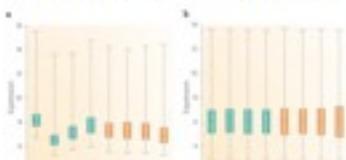
Nature Methods | [Available online](#)

What are Batch Effects



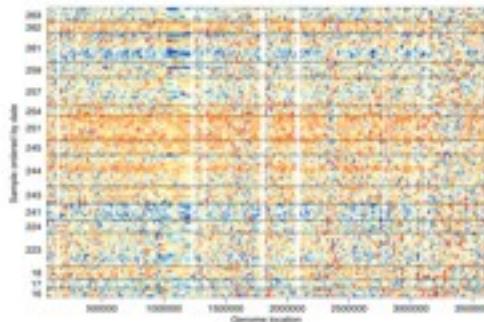
Nature Methods | [Available online](#)

What are Batch Effects



Nature Methods | [Available online](#)

Batch Effects in Sequencing



Math gives us hope

NATURE REVIEWS | GENETICS

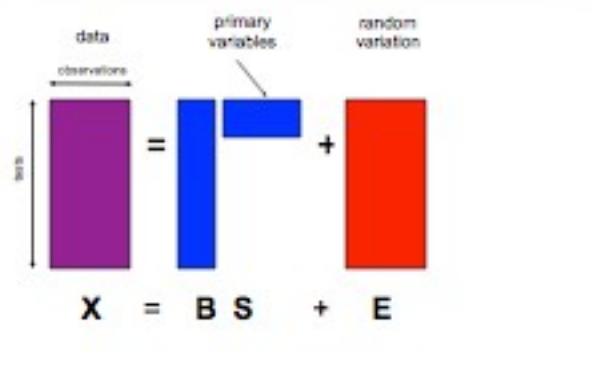
OPINION

Tackling the widespread and critical impact of batch effects in high-throughput data

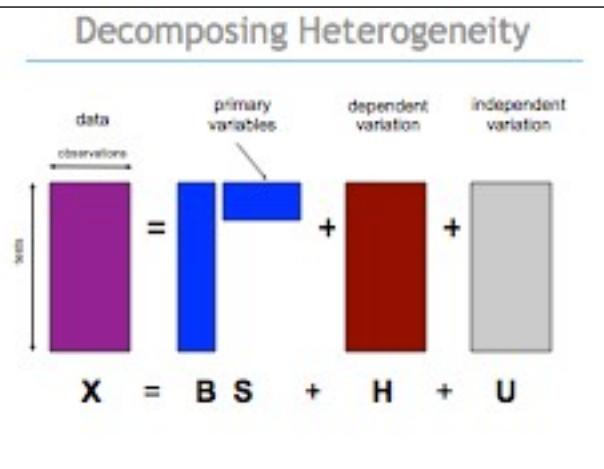
Jeffrey T. Leek, Robert B. Scherf, Héctor Corrado Bravo, David Simcha, Benjamin Langmead, W. Evan Johnson, Donald Geman, Keith Baggerly and Rafael A. Irizarry

Thanks

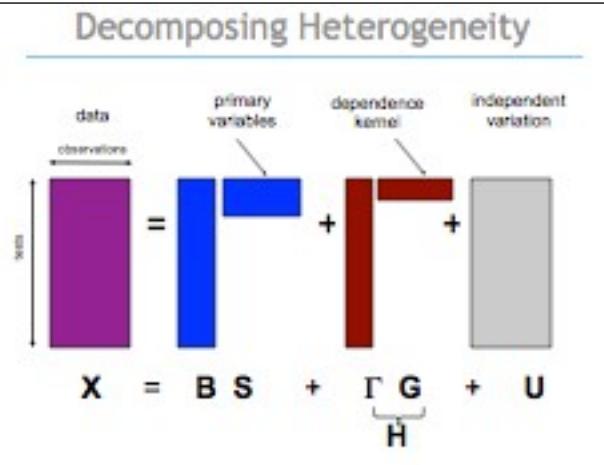
Decomposing Heterogeneity



Decomposing Heterogeneity

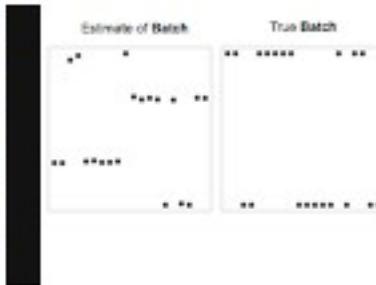
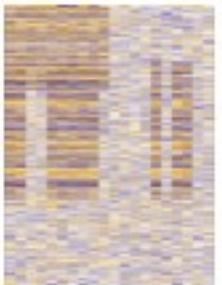


Decomposing Heterogeneity



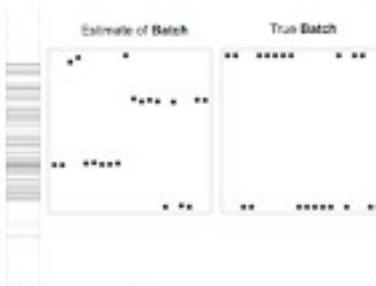
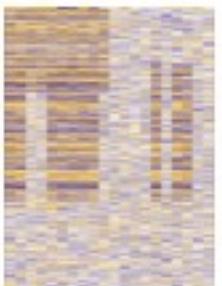
Surrogate Variable Analysis

The Data $P(\text{Group} \& \text{Batch})$



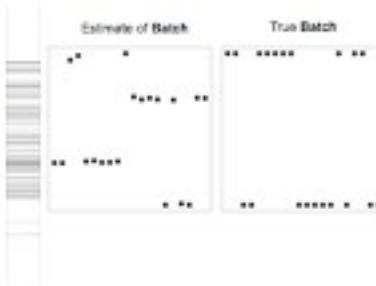
Surrogate Variable Analysis

The Data $P(\text{Group} \& \text{Batch})$



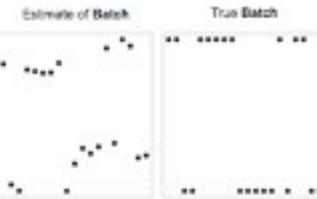
Surrogate Variable Analysis

The Data $P(\text{Group} \& \text{Batch})$



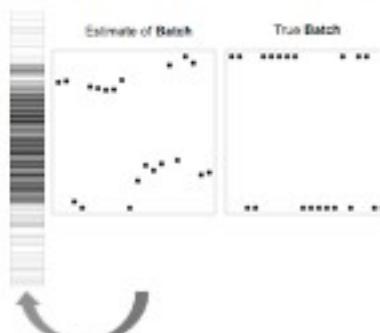
Surrogate Variable Analysis

The Data $P(\text{Group} \& \text{Batch})$



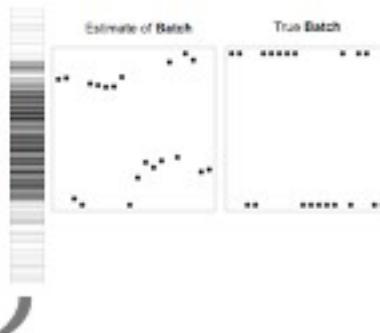
Surrogate Variable Analysis

The Data $P(\text{Group} \& \text{Batch})$



Surrogate Variable Analysis

The Data $P(\text{Group} \& \text{Batch})$

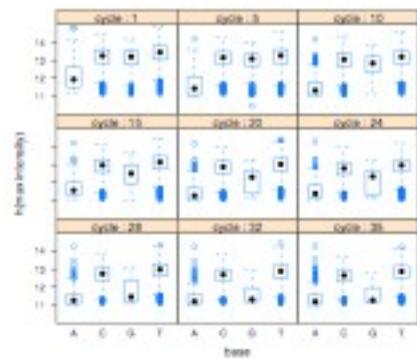


Surrogate Variable Analysis

The Data $P(\text{Group} \mid \text{Batch})$



Base-Cycle Effect



Intensity Model

log intensity read i , cycle j , channel c

$$u_{ijc} = \Delta_{ijc}(\mu_{cja} + x_j^T \alpha_i + \epsilon_{ijc}^\alpha) + (1 - \Delta_{ijc})(\mu_{cjb} + x_j^T \beta_i + \epsilon_{ijc}^\beta)$$

estimate parameters with EM algorithm, which also estimates

$$z_{ijc} := E\{\Delta_{ijc} = 1 | u_{ij}\} = P(\Delta_{ijc} = 1 | u_{ij})$$

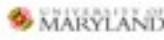
Resources

Papers:

- Corrada Bravo & Irizarry, *Biometrics*, November 2010
- Wu, Irizarry & Corrada Bravo, *Nat. Methods*, May 2010
- Niranjan, et al., *Genome Biology*, Sept. 2011.

Software: <http://cbcb.umd.edu/~hcorrada/seccgen>

- Sfim: Illumina basecalling
- Rsolid: Preprocessing for SOLiD
- Service: Variant calling in targeted pooled samples

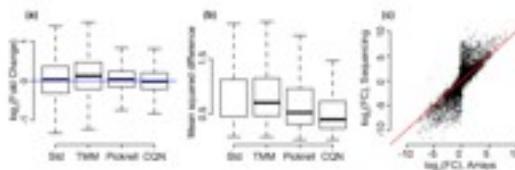


Model

$$Y_{g,i} \mid \mu_{g,i} \sim \text{Poisson}(\mu_{g,i})$$

$$\mu_{g,i} = \exp \left\{ h_i(\theta_{g,i}) + \sum_{j=1}^p f_{i,j}(X_{g,j}) \right\}$$

Results

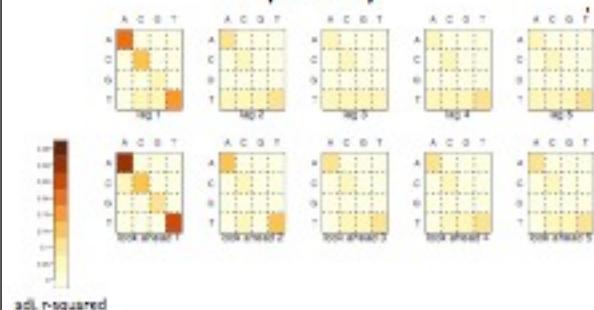


The trick

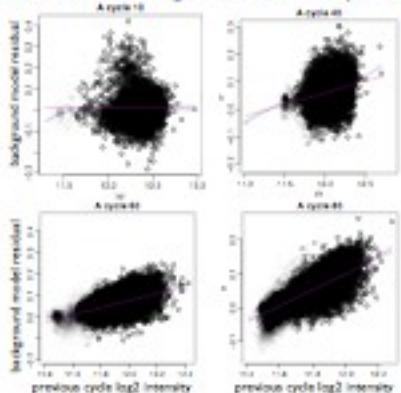
$$\log(Y_{g,i}) \mid \mu_{g,i} \approx \log(\mu_{g,i}) = h_i(\theta_{g,i}) + \sum_{j=1}^p f_{i,j}(X_{g,j}),$$

$$E\left\{\log(Y_{g,i}) - \sum_{j=1}^p f_{i,j}(X_{g,j})\right\} = h_i(\theta_{g,i})$$

Cross-cycle and cross-channel dependency



The next challenge: Residual Dependency



Residual Dependency

