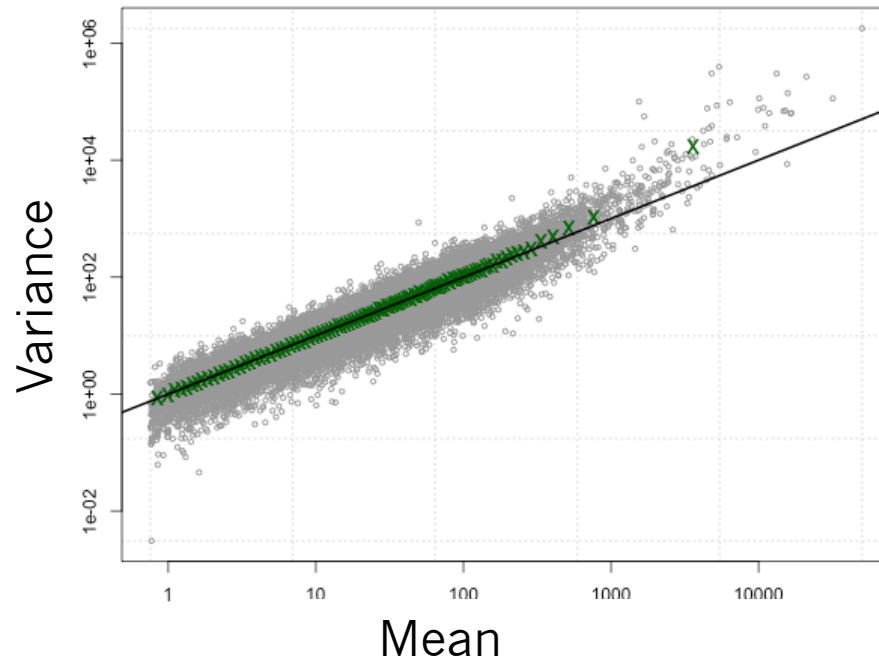# Statistical models for count data analysis (part 2)

- Reminder of tricks used / material already presented:
  - conditional likelihood
  - (local) weighted likelihood
  - linear models
- Bringing them together: a more general framework – GLMs
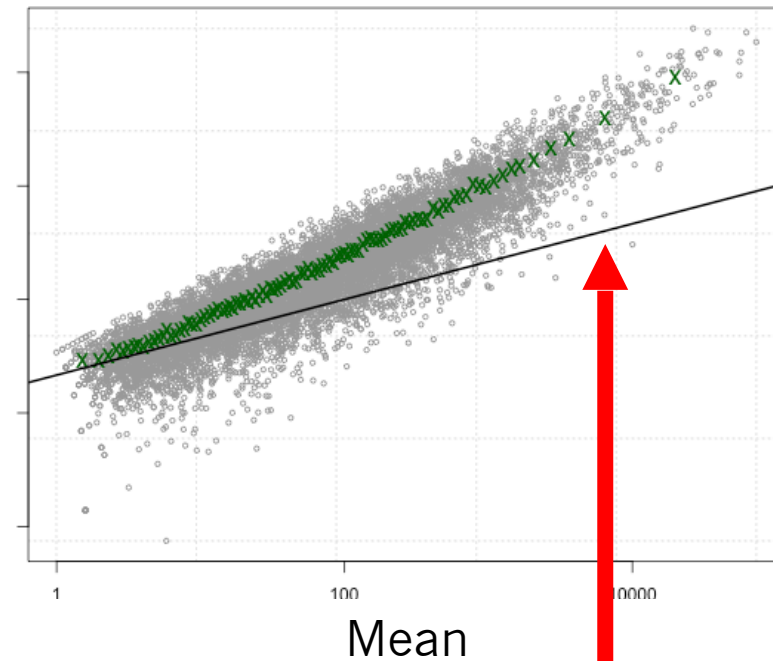- Beyond differential expression: "differential splicing"

# Mean-Variance plots:  What we see in real data



Technical replicates

Biological replicates

mean=variance
(Poisson assumption)

Data from Marioni et al. Genome Research 2008

Data from Parikh et al.
*Genome Biology* 2010

# Conditional likelihood

Likelihood for single **negative binomial** observation:

$$f(y; \mu, \phi) = P(Y = y) = \frac{\Gamma(y + \phi^{-1})}{\Gamma(\phi^{-1})\Gamma(y + 1)} \left( \frac{1}{1 + \mu\phi} \right)^{\phi^{-1}} \left( \frac{\mu}{\phi^{-1} + \mu} \right)^{y}$$

If all libraries are the same size (i.e. $m_i \equiv m$), the sum $Z = Y_1 + \cdots + Y_n \sim \text{NB}(nm\lambda, \phi n^{-1})$

Thus, can form conditional likelihood:

$$l_{Y|Z=z}(\phi) = \left[ \sum_{i=1}^{n} \log \Gamma(y_i + \phi^{-1}) \right] + \log \Gamma(n\phi^{-1}) - \log \Gamma(z + n\phi^{-1}) - n \log \Gamma(\phi^{-1})$$

Statistical Bioinformatics // Institute of Molecular Life Sciences

Log-Likelihood



delta

# Moderated dispersion estimate

Weighted likelihood -- individual log-likelihood plus a weighted version of the common log-likelihood:

Score (1st derivative of LL)
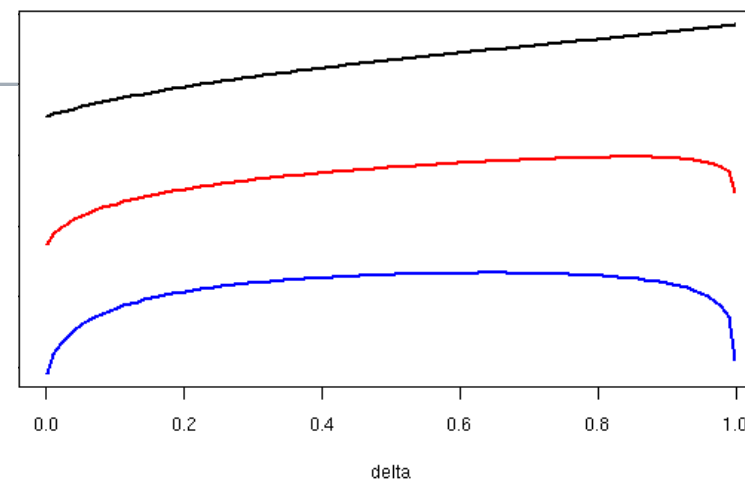


delta

$$\mathrm{WL}(\phi_g) = l_g(\phi_g) + \alpha\, l_C(\phi_g)$$

$l_g$ - quantile-adjusted conditional likelihood

Black: single tag
Blue: common dispersion
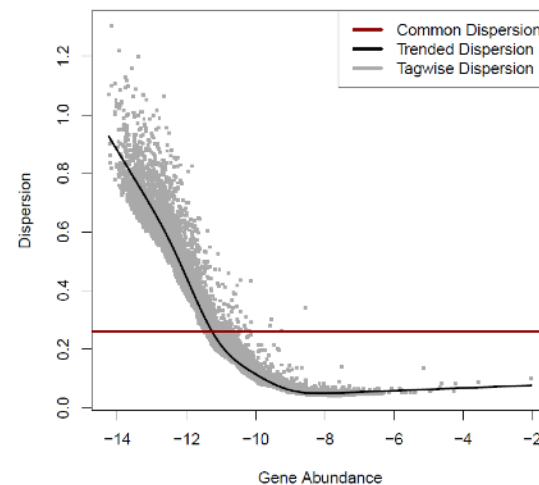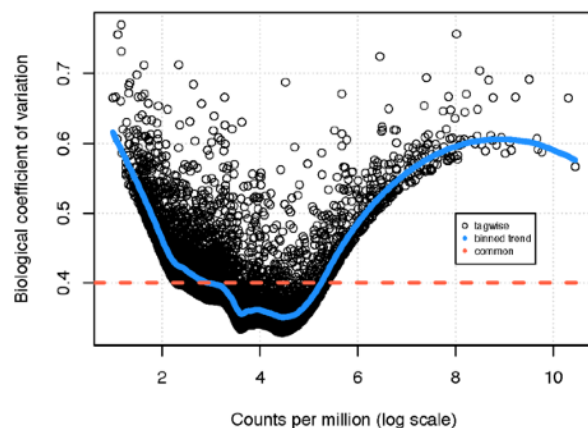Red: Linear combination of the two

$$\delta = \frac{\phi}{\phi+1}$$

# Dispersion varies with mean: moderate (e.g., weighted likelihood) dispersion towards trend

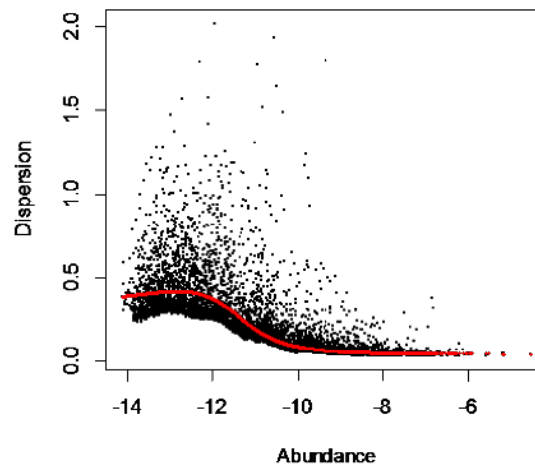$$\text{WL}(\phi_g) = l_g(\phi_g) + \alpha\, l_C(\phi_g)$$

$(1-\alpha)$

Data:
Tuch et al.,
2008



Mouse hemapoeitic
stem cells



Advantage: genes are
allowed to have their
own variance.

Mouse
lymphomas

University of
Zurich^UZH

Statistical Bioinformatics // Institute of Molecular Life Sciences

# Linear Models (microarray setting)

In general, need to specify:

-   Dependent variable

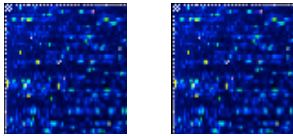-   Explanatory variables (experimental design, covariates, etc.)

More generally:

$$y = X\beta + \epsilon$$

vector of observed data

design matrix

Vector of parameters to estimate

# Analysis of Variance → Linear model

WT x 2                                Cond A x 2                                Cond B x 2



$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \epsilon_4 \\ \epsilon_5 \\ \epsilon_6 \end{bmatrix}$$

$\alpha_1$ = wt log-expression

$\alpha_2$ = Cond A - wt

$\alpha_3$ = Cond B - wt

$E[y_1] = E[y_2] = \alpha_1$          $E[y_3] = E[y_4] = \alpha_1 + \alpha_2$          $E[y_5] = E[y_6] = \alpha_1 + \alpha_3$

Applications: paired designs, multi-factor designs, interactions

—> This particular model only valid for continuous response

# Generalized linear models: a more general framework

Gaussian (normal) distributed response —> various other (common) types.

Three components:

1. Probability distribution of response (in exponential family)

2. Linear predictor (covariates; design matrix)

3. Link function (link mean to linear predictor)

Link function

# Link function and linear predictor

$$E(Y_i) = \mu_i \qquad\qquad g(\mu_i) = \eta_i$$

$$\eta_i = \beta_0 + \beta_1 x_{1i} + ... + \beta_p x_{pi} \qquad\longleftarrow\qquad \text{Linear predictor (covariates)}$$

$$\text{var}(Y_i) = \phi V(\mu)$$

Provides a way to link the mean of response to a linear predictor.

Data is not transformed.

Variance is a function of mean.

# Common distributions, "Canonical" link functions

**Common distributions with typical uses and canonical link functions**

| Distribution | Support of distribution | Typical uses | Link name | Link function | Mean function |
|---|---|---|---|---|---|
| Normal | real: $(-\infty, +\infty)$ | Linear-response data | Identity | $\mathbf{X}\boldsymbol{\beta} = \mu$ | $\mu = \mathbf{X}\boldsymbol{\beta}$ |
| Exponential | real: $(0, +\infty)$ | Exponential-response data, scale parameters | Inverse | $\mathbf{X}\boldsymbol{\beta} = \mu^{-1}$ | $\mu = (\mathbf{X}\boldsymbol{\beta})^{-1}$ |
| Gamma | | | | | |
| Inverse Gaussian | real: $(0, +\infty)$ | | Inverse squared | $\mathbf{X}\boldsymbol{\beta} = \mu^{-2}$ | $\mu = (\mathbf{X}\boldsymbol{\beta})^{-1/2}$ |
| Poisson | integer: $[0, +\infty)$ | count of occurrences in fixed amount of time/space | Log | $\mathbf{X}\boldsymbol{\beta} = \ln(\mu)$ | $\mu = \exp(\mathbf{X}\boldsymbol{\beta})$ |
| Bernoulli | integer: $[0, 1]$ | outcome of single yes/no occurrence | Logit | $\mathbf{X}\boldsymbol{\beta} = \ln\left(\dfrac{\mu}{1-\mu}\right)$ | $\mu = \dfrac{\exp(\mathbf{X}\boldsymbol{\beta})}{1 + \exp(\mathbf{X}\boldsymbol{\beta})} = \dfrac{1}{1 + \exp(-\mathbf{X}\boldsymbol{\beta})}$ |
| Binomial | integer: $[0, N]$ | count of # of "yes" occurrences out of N yes/no occurrences | | | |
| Categorical | integer: $[0, K)$ | outcome of single K-way occurrence | | | |
| | K-vector of integer: $[0, 1]$, where exactly one element in the vector has the value 1 | | | | |
| Multinomial | K-vector of integer: $[0, N]$ | count of occurrences of different types (1 .. K) out of N total K-way occurrences | | | |

http://en.wikipedia.org/wiki/Generalized_linear_model

# RNA-seq setting – Negative binomial regression

Response is negative binomial (dispersion "fixed" to make it in the exponential family).

Link function (relate mean of response to linear combination of parameters)

For example:

$$Y_i \sim \mathrm{NB}(\mu_i, \phi)$$

$$\mathbf{X}\boldsymbol{\beta} = \ln(\mu)$$

X     – design matrix
g()    – link function (here: log)
β     – parameters

**edgeR::glmFit()**

## Same challenge as last time: getting a good estimate of dispersion

Several choices here:

- Maximum Likelihood (MLE)

- Pseudo-Likelihood (PL)

- Quasi-Likelihood (QL)

- Conditional Maximum Likelihood (CML)

- Approximate Conditional Inference (Cox-Reid)

- *quantile-adjusted Maximum Likelihood (qCML)*

$$\mathbf{X}\boldsymbol{\beta} = \ln(\mu)$$

$$Y_i \sim \mathrm{NB}(\mu_i, \phi)$$

$$(\hat{\lambda}_{MLE}, \hat{\phi}_{MLE}) = \arg\max_{\lambda,\phi} l(\lambda, \phi)$$

$$X^2 = \sum_{gij} \frac{(y_{gij} - \hat{\mu}_{gi})^2}{\hat{\mu}_{gi}(1 + \hat{\phi}_{PL}\hat{\mu}_{gi})} = G(n_1 + n_2 - 2)$$

$$D = 2\sum_{gij}\left\{ y_{gij}\log\left[\frac{y_{gij}}{\mu_{gi}}\right] - (y_{gij} + \phi_{QL}^{-1})\log\left[\frac{y_{gij} + \phi_{QL}^{-1}}{\mu_{gi} + \phi_{QL}^{-1}}\right]\right\}$$

# "Cox Reid adjusted profile likelihood" —> Estimation of dispersion parameter

## Parameter Orthogonality and Approximate Conditional Inference

D. R. COX†  and  N. REID

*Imperial College, London*  *University of British Columbia, Vancouver*

[*Read before the* Royal Statistical Society *at a meeting organized by the* Research Section *on Wednesday, 8th October, 1986, Professor A. F. M. Smith in the Chair*]

SUMMARY

We consider inference for a scalar parameter $\psi$ in the presence of one or more nuisance parameters. The nuisance parameters are required to be orthogonal to the parameter of interest, and the construction and interpretation of orthogonalized parameters is discussed in some detail. For purposes of inference we propose a likelihood ratio statistic constructed from the conditional distribution of the observations, given maximum likelihood estimates for the nuisance parameters. We consider to what extent this is preferable to the profile likelihood ratio statistic in which the likelihood function is maximized over the nuisance parameters. There are close connections to the modified profile likelihood of Barndorff-Nielsen (1983). The normal transformation model of Box and Cox (1964) is discussed as an illustration.

*Keywords:* ASYMPTOTIC THEORY; CONDITIONAL INFERENCE; LIKELIHOOD RATIO TEST; NORMAL TRANSFORMATION MODEL; NUISANCE PARAMETERS; ORTHOGONAL PARAMETERS

$$Y_i \sim \text{NB}(\mu_i, \phi)$$

$$\mathbf{X}\beta = \ln(\mu)$$

In this setting, we are trying to get an estimate of dispersion, so the beta (regression) parameters are the "nuisance" parameters.

We turn the problem around later to make inferences about the regression parameters.

$$Y_i \sim \mathrm{NB}(\mu_i, \phi)$$
$$\mathbf{X}\beta = \ln(\mu)$$

# Cox-Reid adjusted profile likelihood

The adjusted profile likelihood (APL) for $\phi_g$ is the penalized log-likelihood

$$\mathrm{APL}_g(\phi_g) = \ell(\phi_g; \mathbf{y}_g, \hat{\beta}_g) - \frac{1}{2}\log \det \mathcal{I}_g.$$

where $\mathbf{y}_g$ is the vector of counts for gene $g$, $\hat{\beta}_g$ is the estimated coefficient vector, $\ell()$ is the log-likelihood function and $\mathcal{I}_g$ is the Fisher information matrix.
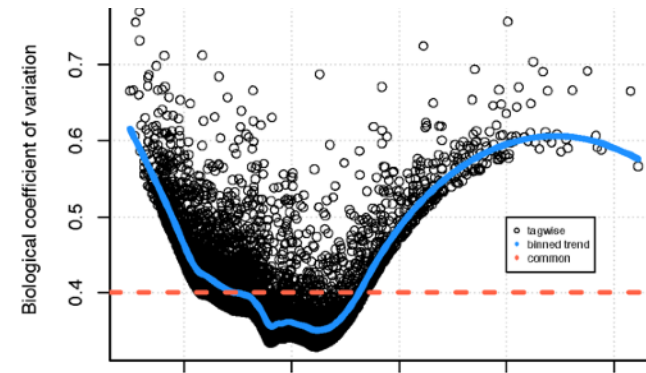
In this approach, $\phi_g$ is estimated by maximizing

$$\mathrm{APL}_g(\phi_g) + G_0 \,\mathrm{APL}_{Sg}(\phi_g),$$

where $G_0$ is the weight given to the shared likelihood and $\mathrm{APL}_{Sg}(\phi_g)$ is the local shared log-likelihood.
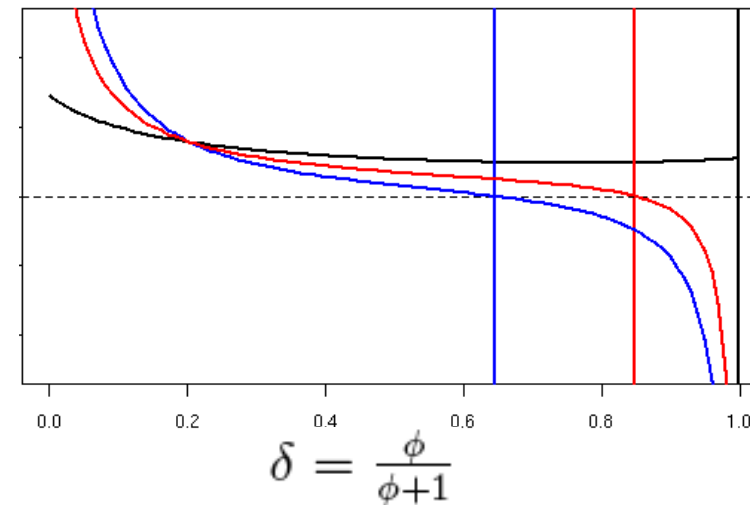
14

# APL is simply another likelihood, so weighted likelihood still works

WL is the individual log-likelihood plus a weighted version of the **common** log-likelihood:

$$\mathrm{WL}(\phi_g) = l_g(\phi_g) + \alpha\, l_C(\phi_g)$$

**$L_g$ - adjusted profile likelihood (or trended version)**

**Black: single tag**
**Blue: common dispersion**
**Red: Linear combination of the two**



Score (1st derivative of LL)



$$\delta = \frac{\phi}{\phi+1}$$

15

# Exponential family

"natural parameter"

$$f(y;\theta) = \exp[a(y)b(\theta) + c(\theta) + d(y)]$$

| Distribution | Natural parameter | $c$ | $d$ |
|---|---|---|---|
| Poisson | $\log\theta$ | $-\theta$ | $-\log y!$ |
| Normal | $\dfrac{\mu}{\sigma^2}$ | $-\dfrac{\mu^2}{2\sigma^2} - \dfrac{1}{2}\log\left(2\pi\sigma^2\right)$ | $-\dfrac{y^2}{2\sigma^2}$ |
| Binomial | $\log\left(\dfrac{\pi}{1-\pi}\right)$ | $n\log\left(1-\pi\right)$ | $\log\dbinom{n}{y}$ |

Optional exercise: what are a(), b() and c() for negative binomial?

Note: negative binomial is NOT in exponential family unless dispersion parameter is treated as fixed.

— from Introduction to Generalized Linear Models, Annette Dobson, 2nd edition.

# Given dispersion estimates (Cox-Reid APL): estimation, statistical testing of regression parameters

Generalized linear model comes with many advantages:

1. Estimation is the same for all response types (so-called Fisher scoring, which effectively turns likelihood maximization into an iteratively re-weighted estimation problem)

2. Asymptotic theory that lead to i) Wald; ii) Score; or, iii) likelihood ratio tests for parameters of interest (more details).  All of these are based on asymptotics ("large" sample approximations) – how to choose one that works well in practice?

## Large sample theory – Result 1 (Regression parameter estimates are asymptotically normal)

The Wald test follows immediately from the fact that the information matrix for generalized linear models is given by

$$\mathbf{I}(\boldsymbol{\beta}) = \mathbf{X}'\mathbf{W}\mathbf{X}/\phi, \qquad\qquad (B.9)$$

so the large sample distribution of the maximum likelihood estimator $\hat{\boldsymbol{\beta}}$ is multivariate normal

$$\hat{\boldsymbol{\beta}} \sim N_p(\boldsymbol{\beta}, (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\phi). \qquad\qquad (B.10)$$

with mean $\boldsymbol{\beta}$ and variance-covariance matrix $(\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\phi$.

Tests for subsets of $\boldsymbol{\beta}$ are based on the corresponding marginal normal distributions.

(Wald test used in DESeq2 package)

$$\mathcal{I}(\theta) = \mathbb{E}\left\{\left.\left[\frac{\partial}{\partial\theta}\log L(\theta;X)\right]^2\right|\theta\right\}.$$

Statistical Bioinformatics // Institute of Molecular Life Sciences

# Large sample theory – Result 2 (score is asymptotically normal)

$$\dot{\ell}_1 = \frac{\partial\ell}{\partial\boldsymbol{\theta}_1}$$

The "score" function is the first derivative (gradient) of the log-likelihood function, is (asymptotically) normally distributed with mean 0 and variance(-covariance) Fisher information.

$$\dot{\ell}_2 = \frac{\partial\ell}{\partial\boldsymbol{\theta}_2}$$

Say, we to test $H_0$: $\theta_2 = 0$, $\theta_1$ is/are "nuisance" parameter(s)

$$\mathcal{I}_{2.1} = \mathcal{I}_{22} - \mathcal{I}_{21}\mathcal{I}_{11}^{-1}\mathcal{I}_{12}.$$

$$\mathcal{I} = \begin{pmatrix} \mathcal{I}_{11} & \mathcal{I}_{12} \\ \mathcal{I}_{21} & \mathcal{I}_{22} \end{pmatrix}$$

$$S = \dot{\ell}_2^T \mathcal{I}_{2.1}^{-1} \dot{\ell}_2$$

# Large sample theory – Result 3 (likelihood ratio test)

$$D = -2 \ln \left( \frac{\text{likelihood for null model}}{\text{likelihood for alternative model}} \right)$$
$$= -2 \ln(\text{likelihood for null model}) + 2 \ln(\text{likelihood for alternative model})$$

http://en.wikipedia.org/wiki/Likelihood-ratio_test

General form (exponential family)

$$-2 \log \lambda = 2 \sum_{i=1}^{n} \frac{y_i(\tilde{\theta}_i - \hat{\theta}_i) - b(\tilde{\theta}_i) + b(\hat{\theta}_i)}{a_i(\phi)}$$

**edgeR::glmLRT()**

Again, large sample theory says this is approx. $\chi^2$ with degrees of freedom according to the difference in the number of parameters between null and alternative (assuming they are nested).

# Some interesting generalizations of NB modeling for RNA-seq (1)

If $Y_{ijk}$ has a Poisson distribution, then $\mathrm{Var}(Y_{ijk}) = \mu_{ik}$.

If $Y_{ijk}$ has an NB2 distribution, then $\mathrm{Var}(Y_{ijk}) = \mu_{ik}(1 + \phi\mu_{ik})$.

if $Y_{ijk}$ has an NBP distribution, then $\mathrm{Var}(Y_{ijk}) = \mu_{ik}(1 + \phi\mu_{ik}^{\alpha-1})$.

(generalization of the model:
mean-variance relationship)

Di et al., SAGMB 2011 10(1): 24

# Some interesting generalizations of NB modeling for RNA-seq (2)

$$\lambda = 2(l(\hat{\beta}) - l(\tilde{\beta})),$$

$$r = \text{sign}(\hat{\psi} - \psi_0)\sqrt{\lambda}$$

## Higher order asymptotics

For testing a one-dimensional parameter of interest ($q = 1$), Barndorff-Nielsen (1986, 1991) showed that a *modified directed deviance*

$$r^* = r - \frac{1}{r}\log\left(z\right) \tag{5}$$

is, in wide generality, asymptotically standard normally distributed to a higher order of accuracy than the directed deviance $r$ itself, where $z$ is an adjustment term to be discussed below. Tests based on high-order asymptotic adjustment to the likelihood ratio statistic, such as $r^*$ or its approximation (explained below), are referred to as higher-order asymptotic (HOA) tests. They generally have better accuracy than corresponding unadjusted likelihood ratio tests, especially in situations where the sample size is small and/or when the number of nuisance parameters ($p–q$) is large.

Di et al., SAGMB 2013; 12(1): 49–70

# Some interesting generalizations of NB modeling for RNA-seq (3)

$$LRT_k = 2\left(\ell_k(\hat{\mu}_k|y_k) - \ell_k(\tilde{\mu}_k|y_k)\right) \quad \longrightarrow \quad LRT_k \sim \Phi_k \chi_q^2 + O_p(n^{-1/2})$$

$$\hat{\Phi}_k = \frac{2\left(\ell_k(y_k|y_k) - \ell_k(\hat{\mu}_k|y_k)\right)}{n-p}$$

Accounting for the uncertainty in estimating dispersion

$$F_{QL} = \frac{LRT_k/q}{\hat{\Phi}_k}$$

Lund et al., SAGMB 2012; 11(5):8

**University of Zurich**[UZH]

## Offset ξ explicitly in GLM:

$$E[Y] = \mu = g^{-1}(\eta) = g^{-1}(X\beta + \xi)$$

# Some interesting generalizations of NB modeling for RNA-seq (4)



Integrate sample-specific normalization via offset

Profiles vary from sample to sample:
GC content
Gene length

DOES NOT change data, use offsets to modify expected mean

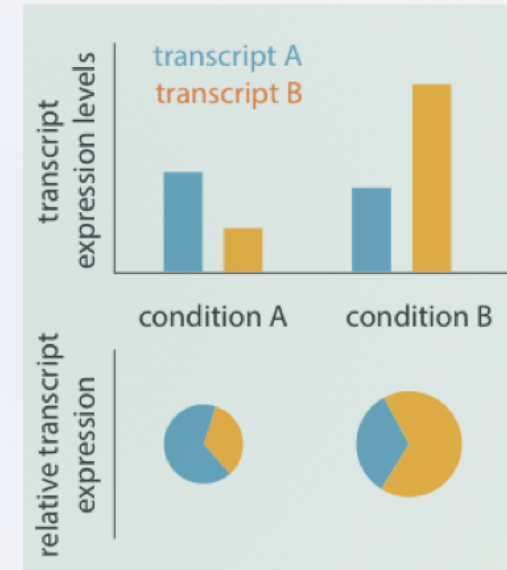Give a sample (or gene)-specific offset to edgeR/DESeq2

Hansen, Irizarry, Wu Biostatistics 2012          24

# Some terms: DTE, DEU, DTU



**Differential transcript expression (DTE)**

**Differential exon usage (DEU)**

**Differential transcript usage (DTU)**
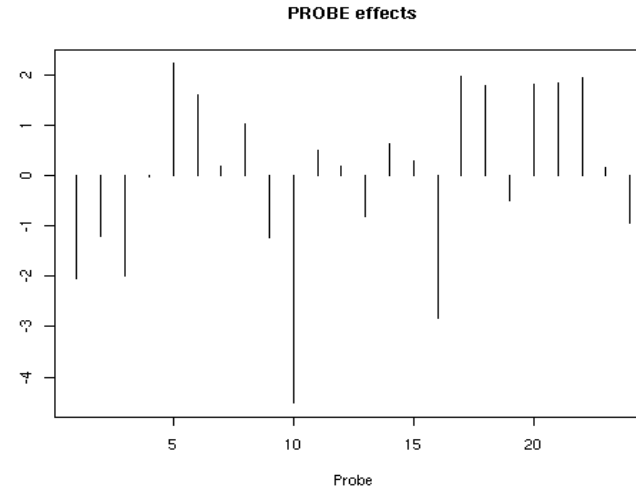
**differential splicing**

Statistical Bioinformatics // Institute of Molecular Life Sciences
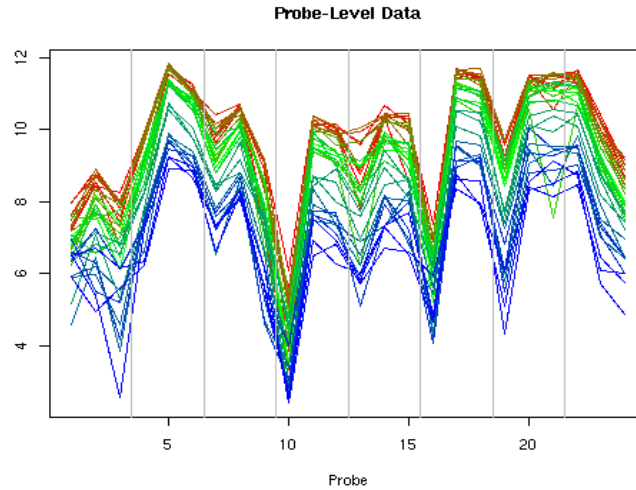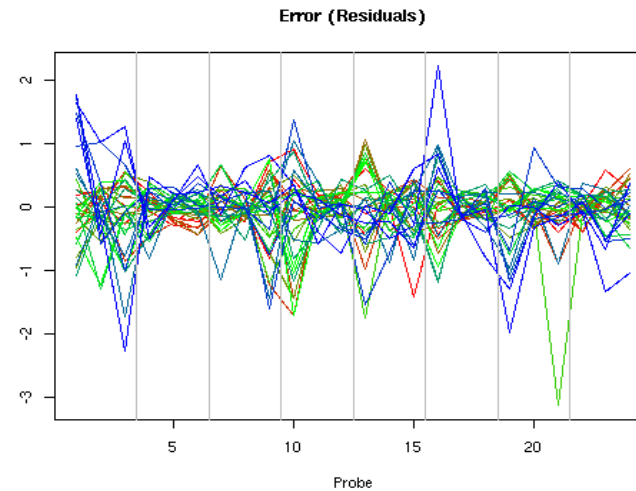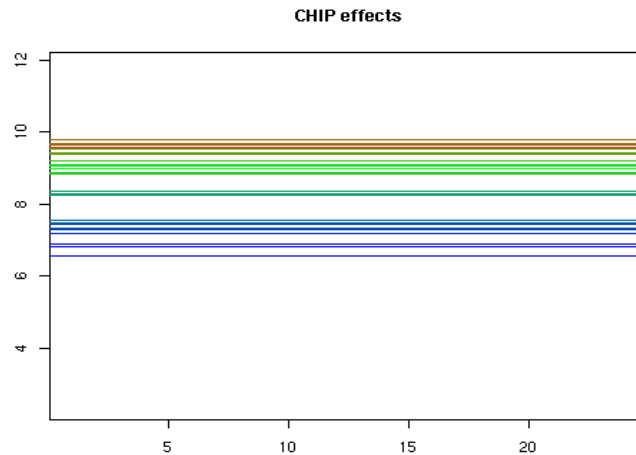


- Data for gene that is DE between heart (red=100% heart) and brain (blue=100% brain).
- 11 mixtures x 3 replicates = 33 samples (33 lines)
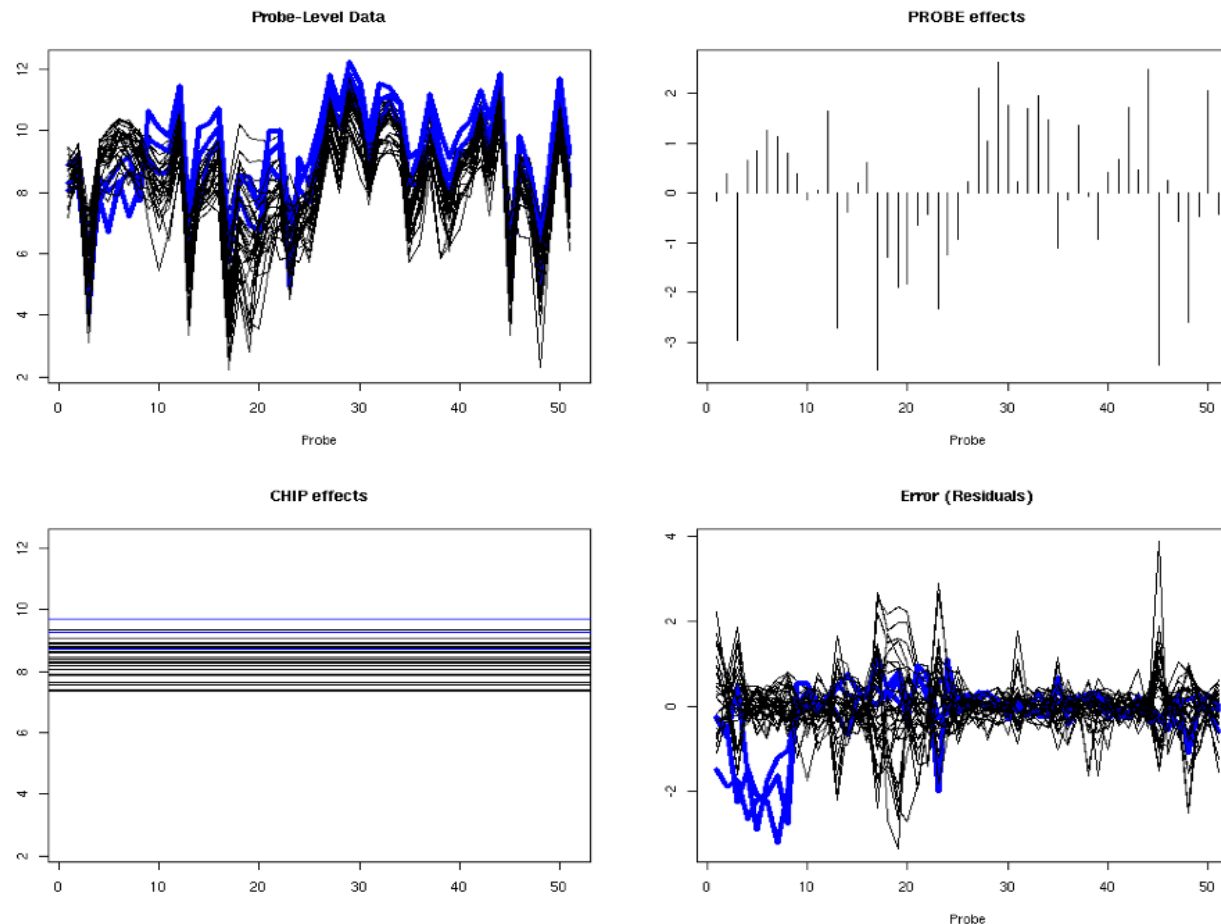- Note the parallelism: probes have different affinities

Tissue mixture dataset

26

# (Digression 2/3) Differential expression: Affy microarrays



$$y_{ik} = g_i + p_k + e_{ik}$$

# Digression 3/3: "Differential splicing" or "Differential isoform usage": Affy microarrays



$$y_{ik} = g_i + p_k + e_{ik}$$

# (back to RNA-seq) Beyond differential expression: differential splicing



**Prediction of alternative isoforms from exon expression levels in RNA-Seq experiments**

Hugues Richard[1,*], Marcel H. Schulz[1,2], Marc Sultan[3], Asja Nürnberger[3], Sabine Schrinner[3], Daniela Balzereit[3], Emilie Dagand[3], Axel Rasche[3], Hans Lehrach[3], Martin Vingron[1], Stefan A. Haas[1] and Marie-Laure Yaspo[3]

[1]Department of Computational Molecular Biology, Max Planck Institute for Molecular Genetics, Ihnestr. 73,
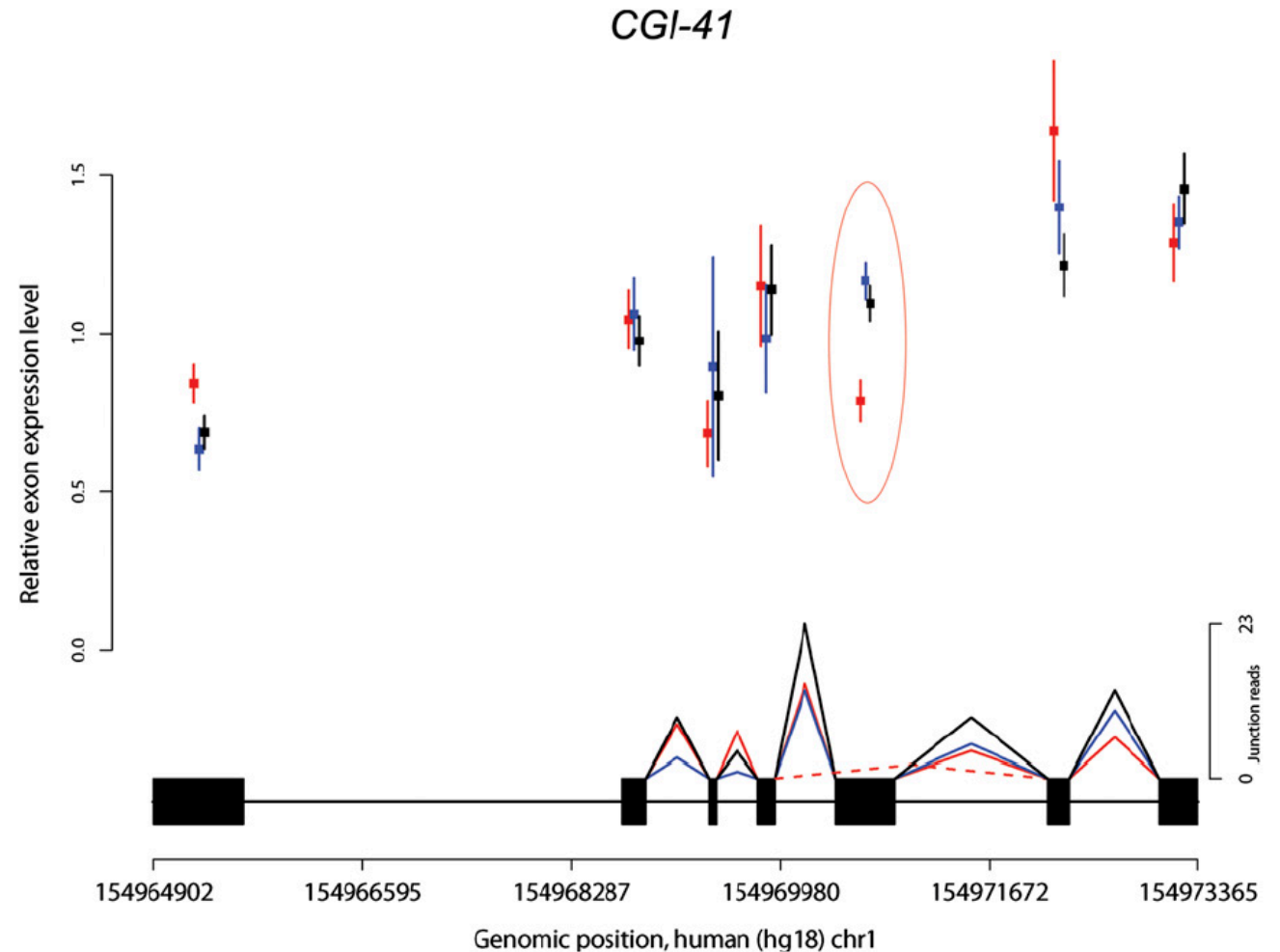[2]International Max Planck Research School for Computational Biology and Scientific Computing, and
[3]Department of Vertebrate Genomics, Max Planck Institute for Molecular Genetics, Ihnestr. 73, 14195 Berlin, Germany

**Sex-specific and lineage-specific alternative splicing in primates**

Ran Blekhman,[1,4,5] John C. Marioni,[1,4,5] Paul Zumbo,[2] Matthew Stephens,[1,3,5] and Yoav Gilad[1,5]

[1]Department of Human Genetics, University of Chicago, Chicago, Illinois 60637, USA; [2]Keck Biotechnology Laboratory, New Haven, Connecticut 06511, USA; [3]Department of Statistics, University of Chicago, Chicago, Illinois 60637, USA

# Counting: a few considerations (exon-level)

All the downstream statistical methods start with a count table.

How to get one?

- annotation-based? What about novel genes?
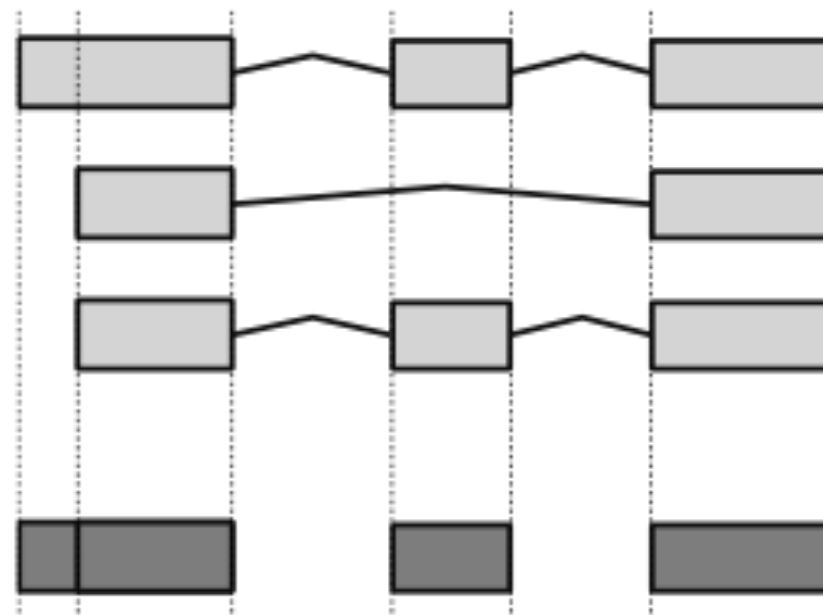- gene-level versus transcript-level? versus exon-level?
- ambiguities
- junctions?



**Figure 1.** Flattening of gene models: This (fictional) gene has three annotated transcripts involving three exons (light shading), one of which has alternative boundaries. We form counting bins (dark shaded boxes) from the exons as depicted; the exon of variable length gets split into two bins.

Anders et al. 2012 Genome Research

## DEXSeq

## Transcript inventory versus differential expression

Shotgun RNA-seq data can be used both for identification of transcripts and for differential expression analysis. In the former, one annotates the regions of the genome that can be expressed, i.e., the exons, and how the pre-mRNAs are spliced into transcripts. In differential expression analysis, one aims to study the regulation of these processes across different conditions. For the method described here, we assume that a transcript inventory has already been defined, and focus on differential expression.

# DEXSeq – general structure: exon-level models

We use generalized linear models (GLMs) (McCullagh and Nelder 1989) to model read counts. Specifically, we assume $K_{ijl}$ to follow a negative binomial (NB) distribution:

$$K_{ijl} \sim NB\left(\text{mean} = s_j \mu_{ijl}, \text{dispersion} = \alpha_{il}\right), \tag{1}$$

where $\alpha_{il}$ is the dispersion parameter (a measure of the distribution's spread; see below) for counting bin $(i, l)$, and the mean is predicted via a log-linear model as

$$\log \mu_{ijl} = \beta_i^G + \beta_{il}^E + \beta_{i\rho_j}^C + \beta_{i\rho_j l}^{EC}. \tag{2}$$

$i$ – gene
$j$ – sample … $\rho_j$ is condition (categorical)
$l$ – bin

$\beta^G$ – baseline "expression strength"
$\beta^E$ – "exon" (bin) effect
$\beta^C$ – condition effect
$\beta^{EC}$ – condition x "exon" interaction

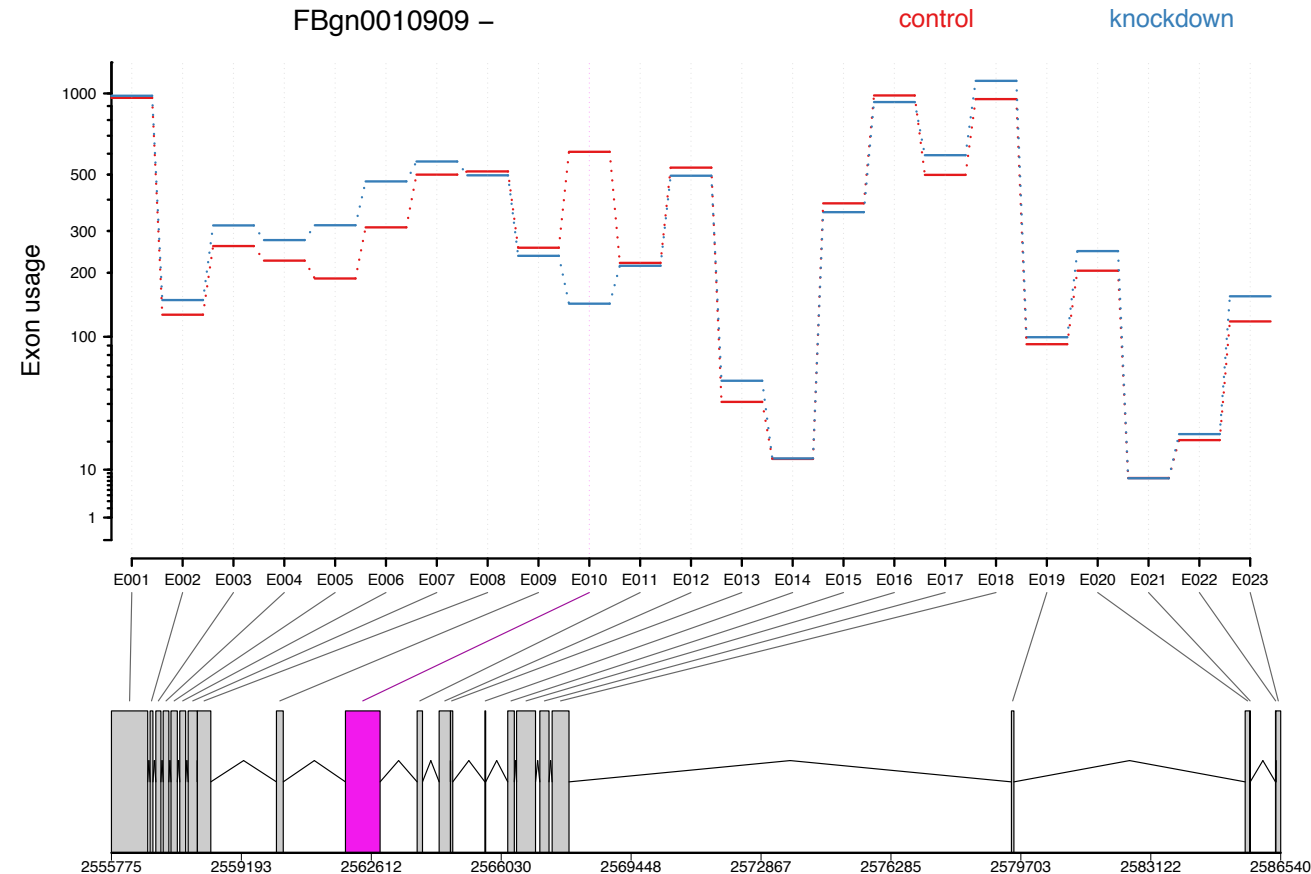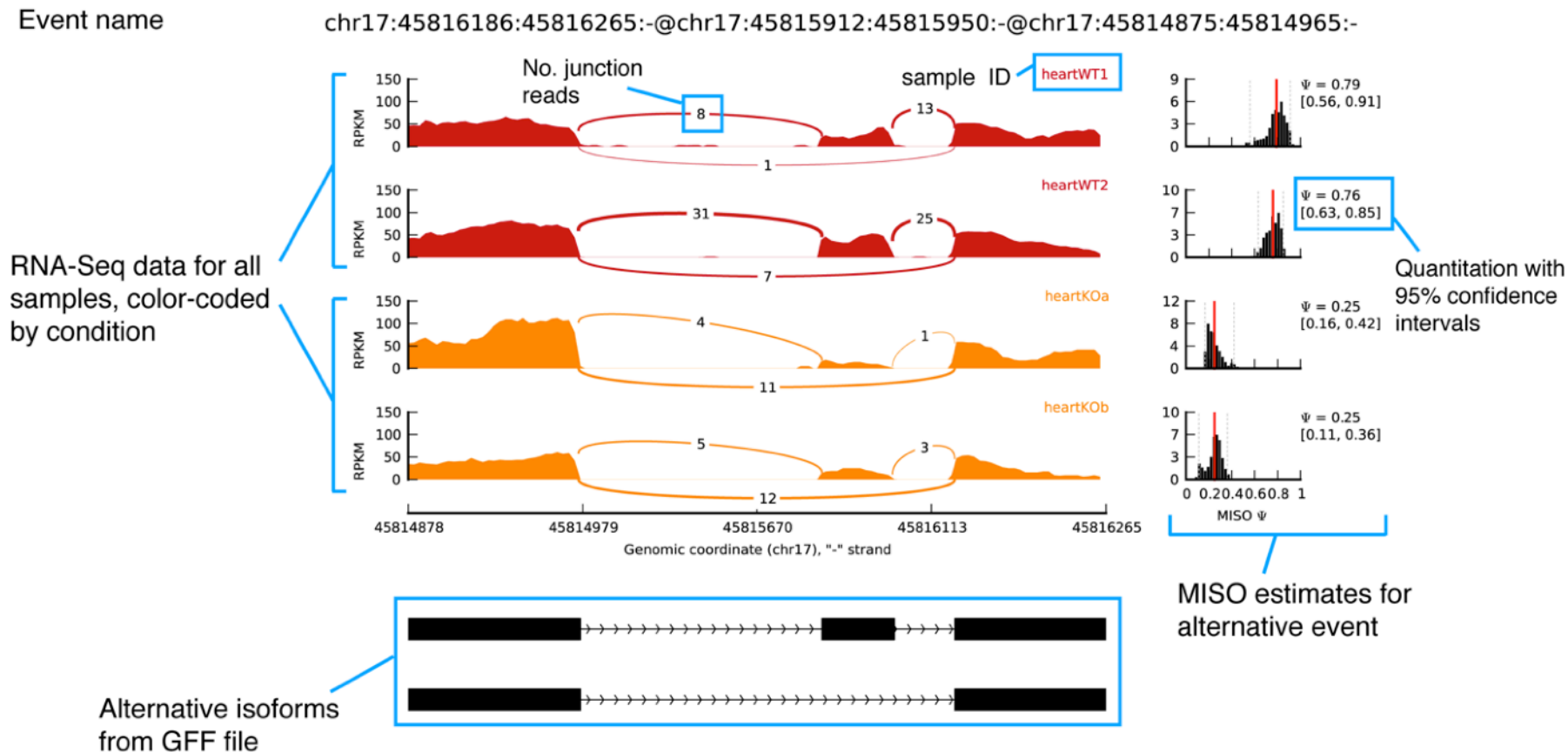# DEXSeq: sig. interaction terms = differential exon usage

(DEXSeq vignette)



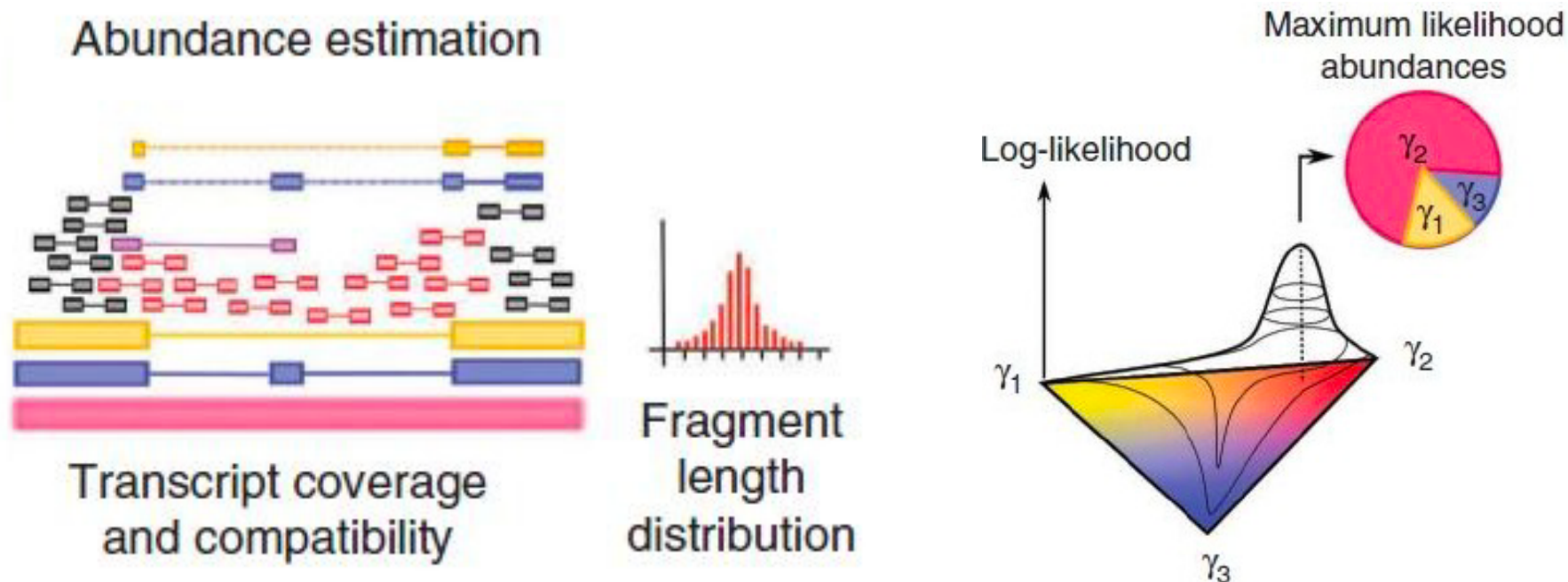**Figure 6: Fitted splicing**
The plot represents the estimated effects, as in Figure 3, but after subtraction of overall changes in gene expression.

33

# Percent spliced in (psi) -- MISO
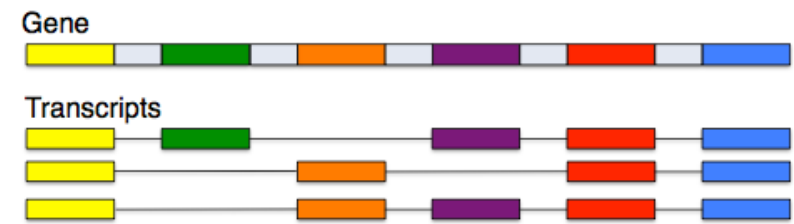


http://genes.mit.edu/burgelab/miso/docs/: "currently, MISO does not handle replicates / groups of samples in any special way" —> rMATs (Shen et al., PNAS, 2014)

# Isoform-level estimation: cufflinks (kallisto, salmon, RSEM), cuffdiff2; many others



Abundance estimation

Transcript coverage and compatibility

Fragment length distribution

Log-likelihood

Maximum likelihood abundances

From estimated isoform abundance from set of (assembled) transcripts, use Jenson-Shannon (JS) divergence to determine change in the mix of transcripts between conditions.

# DTU —> dirichlet-multinomial distribution

Estimated:

- transcript ratios

$$\Pi = (\pi_1, \pi_2, \pi_3)$$

Observed:

- transcript counts $\qquad Y = (y_1, y_2, y_3)$

- gene expression $\qquad n = \sum_{j=1}^{k} y_j$

Multinomial:
$$P(\mathbf{Y} = \mathbf{y} | \mathbf{\Pi} = \pi) = \binom{n}{\mathbf{y}} \prod_{j=1}^{k} \pi_j^{y_j}$$

Dirichlet:
$$P(\mathbf{\Pi} = \pi) = \frac{\Gamma(\gamma_+)}{\prod_{j=1}^{k} \Gamma(\gamma_j)} \prod_{j=1}^{k} \pi_j^{\gamma_j - 1}, \gamma_+ = \sum_{j=1}^{k} \gamma_j$$

Dirichlet-multinomial:
$$P(\mathbf{Y} = \mathbf{y}) = \binom{n}{\mathbf{y}} \frac{\Gamma(\gamma_+)}{\Gamma(n + \gamma_+)} \prod_{j=1}^{k} \frac{\Gamma(y_j + \gamma_j)}{\Gamma(\gamma_j)}, \gamma_j = \pi_j \gamma_+$$